

Dimensionality Reduction: Singular Value Decomposition (SVD) & Principal Component Analysis (PCA)

High-Dimensional Data

High-Dimensions = Lot of Features

- Document classification
 - Features per document = thousands of words/unigrams, millions of bigrams, contextual information
- Surveys – Netflix
 - 480189 users x 17770 movies

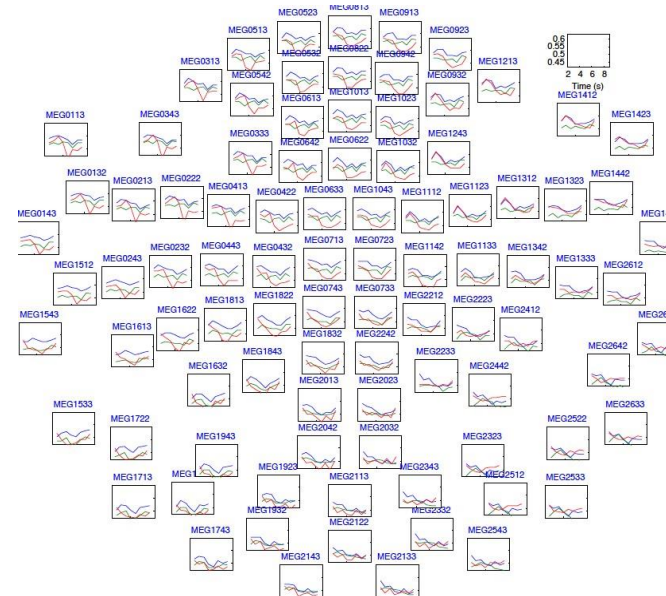
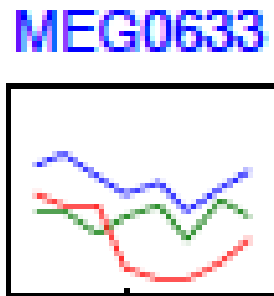
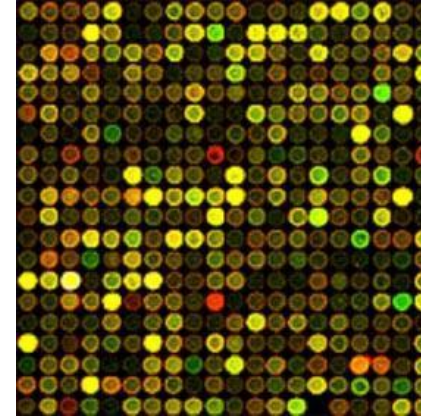


	movie 1	movie 2	movie 3	movie 4	movie 5	movie 6
Tom	5	?	?	1	3	?
George	?	?	3	1	2	5
Susan	4	3	1	?	5	1
Beth	4	3	?	2	4	2

High-Dimensional Data

High-Dimensions = Lot of Features

- Discovering gene networks
 - 10,000 genes x 1000 drugs x several species
- MEG Brain Imaging
 - 120 locations x 500 time points x 20 objects



Curse of Dimensionality

Why are more features bad?

- Redundant features
 - not all words are useful to classify a document
- Hard to interpret and visualize
- Hard to store and process data
 - computationally challenging
- Complexity of decision boundaries tends to grow with # features

Dimensionality Reduction

Represent data with fewer dimensions

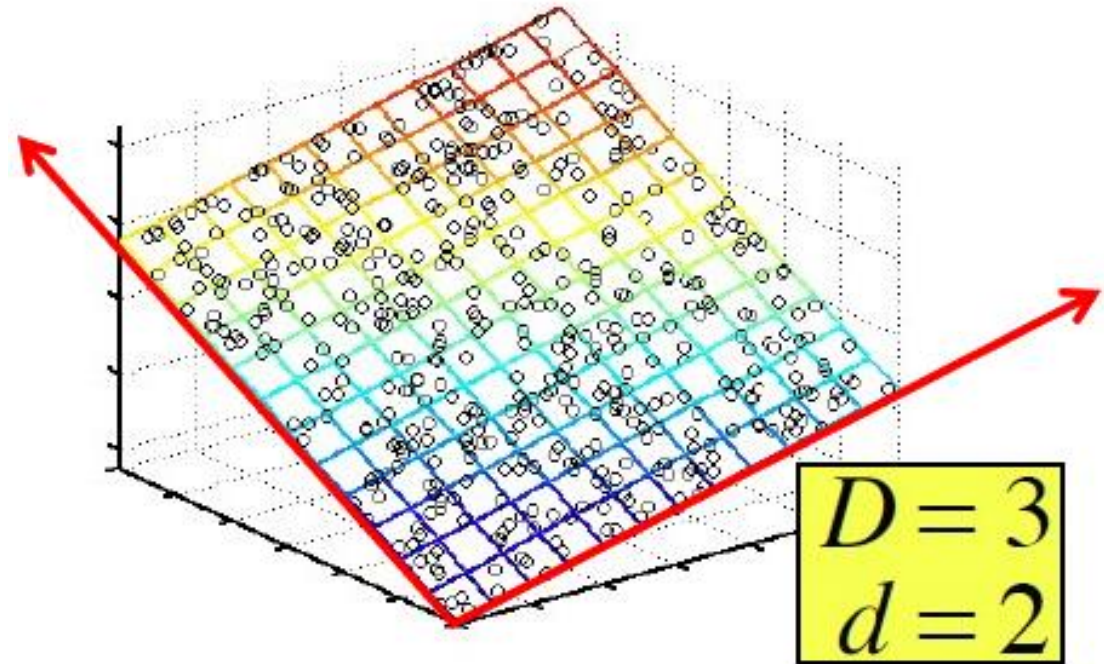
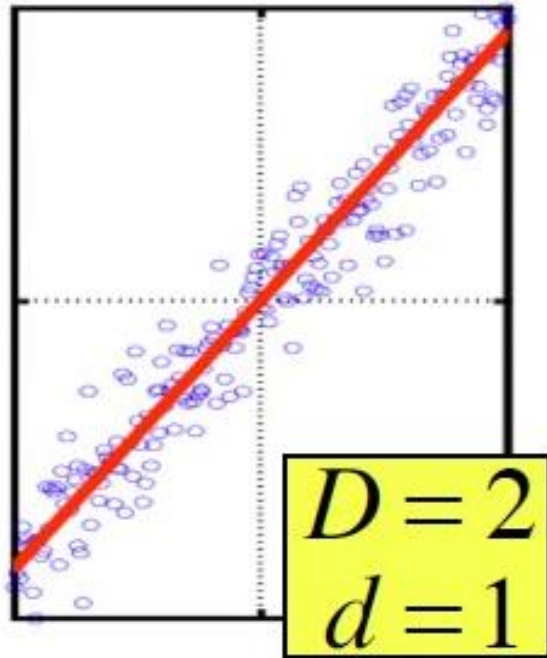
- **Easier learning** – fewer parameters
- **Visualization** – show high dimensional data in 2D
- Discover “**intrinsic dimensionality**” of data
 - high dimensional data that is truly lowerdimensional
 - noise reduction

Informally: Given data points in D -dimensional space

- Convert them to data points in $d < D$ dimensions
- Goal: With minimal/maximal information loss/maintenance

Dimensionality Reduction

Assumption: Data (approximately) lies on a lower dimensional space



Singular Value Decomposition(SVD)

Orthonormal Basis

A set of vectors $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d\} \subset \mathbb{R}^d$ form an orthonormal basis if:

- unit ℓ_2 -norm: $\|\mathbf{u}_i\|_2 = 1$, for all $i=1$ to d
- orthogonal to each other: $(\mathbf{u}_i, \mathbf{u}_j) = 0$, for all $i \neq j$

Let $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d\}$ be any orthonormal basis of \mathbb{R}^d

- Any vector $\mathbf{x} \in \mathbb{R}^d$ can be expressed as a linear combination of $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d$, for some real-valued $\alpha_1, \alpha_2, \dots, \alpha_d$

$$\mathbf{x} = \alpha_1 \mathbf{u}_1 + \alpha_2 \mathbf{u}_2 + \dots + \alpha_d \mathbf{u}_d$$

Singular Value Decomposition (SVD)

- Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be any matrix
- Rank: $r = \text{rank}(\mathbf{A}) (\leq m, n)$
- SVD: $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$
 - Singular values: $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$
 - $\mathbf{\Sigma}$ diagonal matrix, with $\Sigma_{ii} = \sigma_i$
 - Left singular matrix/vectors: $\mathbf{U} = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\} \in \mathbb{R}^m$ forms an orthonormal basis, i.e., $\mathbf{U}^T \mathbf{U} = \mathbf{I}_r$
 - Right singular matrix/vectors: $\mathbf{V} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\} \in \mathbb{R}^n$ forms an orthonormal basis, i.e., $\mathbf{V}^T \mathbf{V} = \mathbf{I}_r$

Singular Value Decomposition (SVD)

- Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be any matrix
- Rank: $r = \text{rank}(\mathbf{A}) (\leq m, n)$
- SVD: $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$
 - Singular values: $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$

The diagram illustrates the SVD decomposition of matrix \mathbf{A} into a sum of rank-1 matrices. It shows the equation:

$$\mathbf{A}_{m \times n} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \dots$$

Each term in the sum is represented by a brown rectangular box containing the expression $\sigma_i \mathbf{u}_i \mathbf{v}_i^T$. The boxes are arranged horizontally, separated by plus signs. An orange arrow points to the right below the boxes, labeled "decreasing importance", indicating that the singular values σ_i decrease in magnitude as i increases, and thus the corresponding rank-1 components become less important in approximating the matrix \mathbf{A} .

Singular Value Decomposition (SVD)

- Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be any matrix
- Rank: $r = \text{rank}(\mathbf{A}) (\leq m, n)$
- SVD: $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$
- Truncated SVD: for $0 < k < r$
 - $\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$
 - \mathbf{A}_k is the **best** rank- k approximation to \mathbf{A}

$$\mathbf{A}_k = \arg \min_{\mathbf{B}} \|\mathbf{A} - \mathbf{B}\|_F^2, \text{ s.t. } \text{rank}(\mathbf{B}) \leq k.$$

Matrix Frobenius Norm

- Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be any matrix

- Frobenius norm

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2}$$

A generalization of the l_2 -vector norm to matrix

- Relation with singular values

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^r \sigma_i^2}$$

σ_i is the i -th singular value of \mathbf{A}

Error of Truncated SVD

- Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be any matrix and $\text{rank}(\mathbf{A})=r$
- SVD: $\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$
- Truncated SVD: $\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ (for $0 < k < r$)
- Error of the truncated SVD

$$\|\mathbf{A} - \mathbf{A}_k\|_F^2 = \sum_{i=k+1}^r \sigma_i^2$$

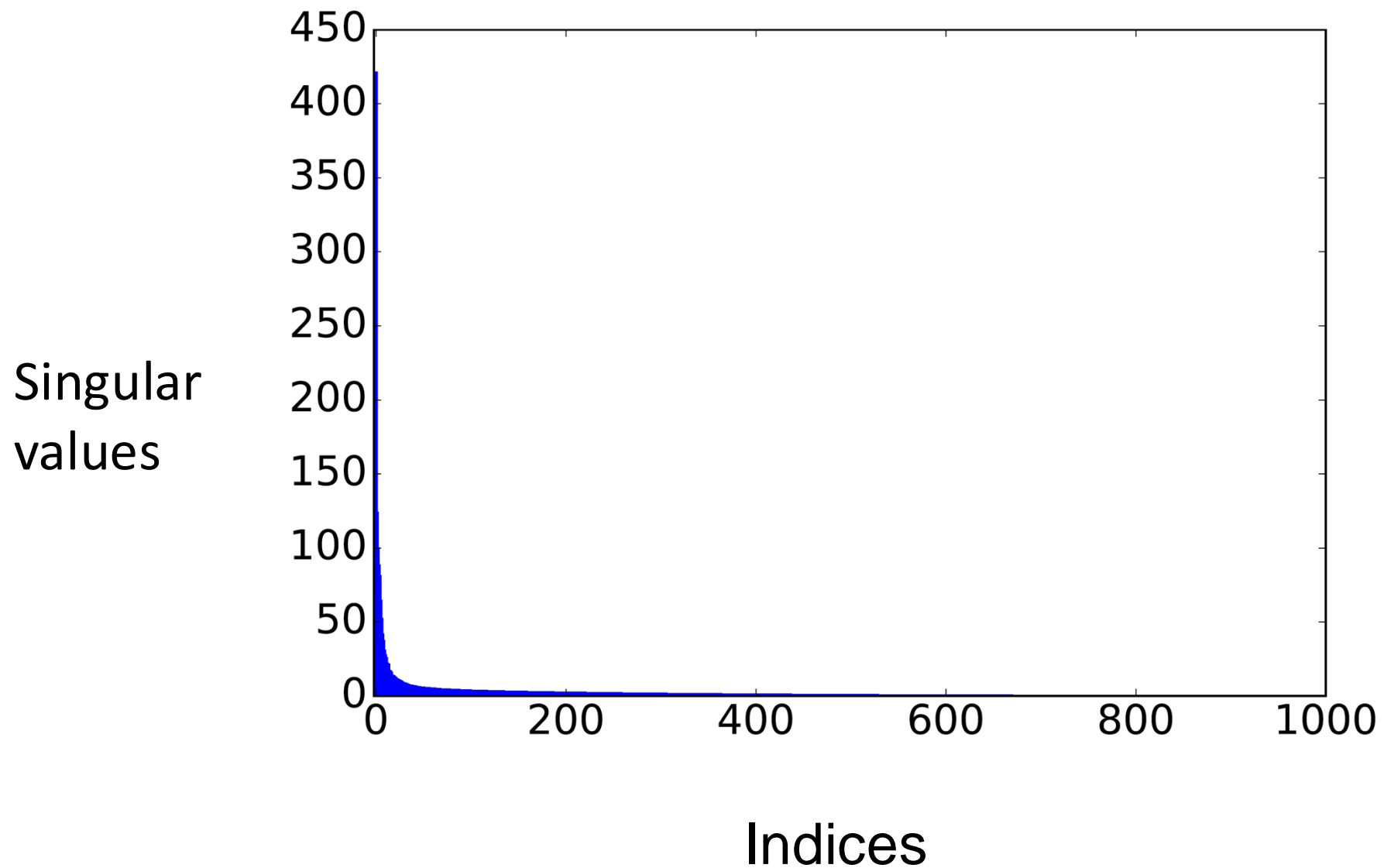
Bottom (the smallest) singular values

SVD: Example

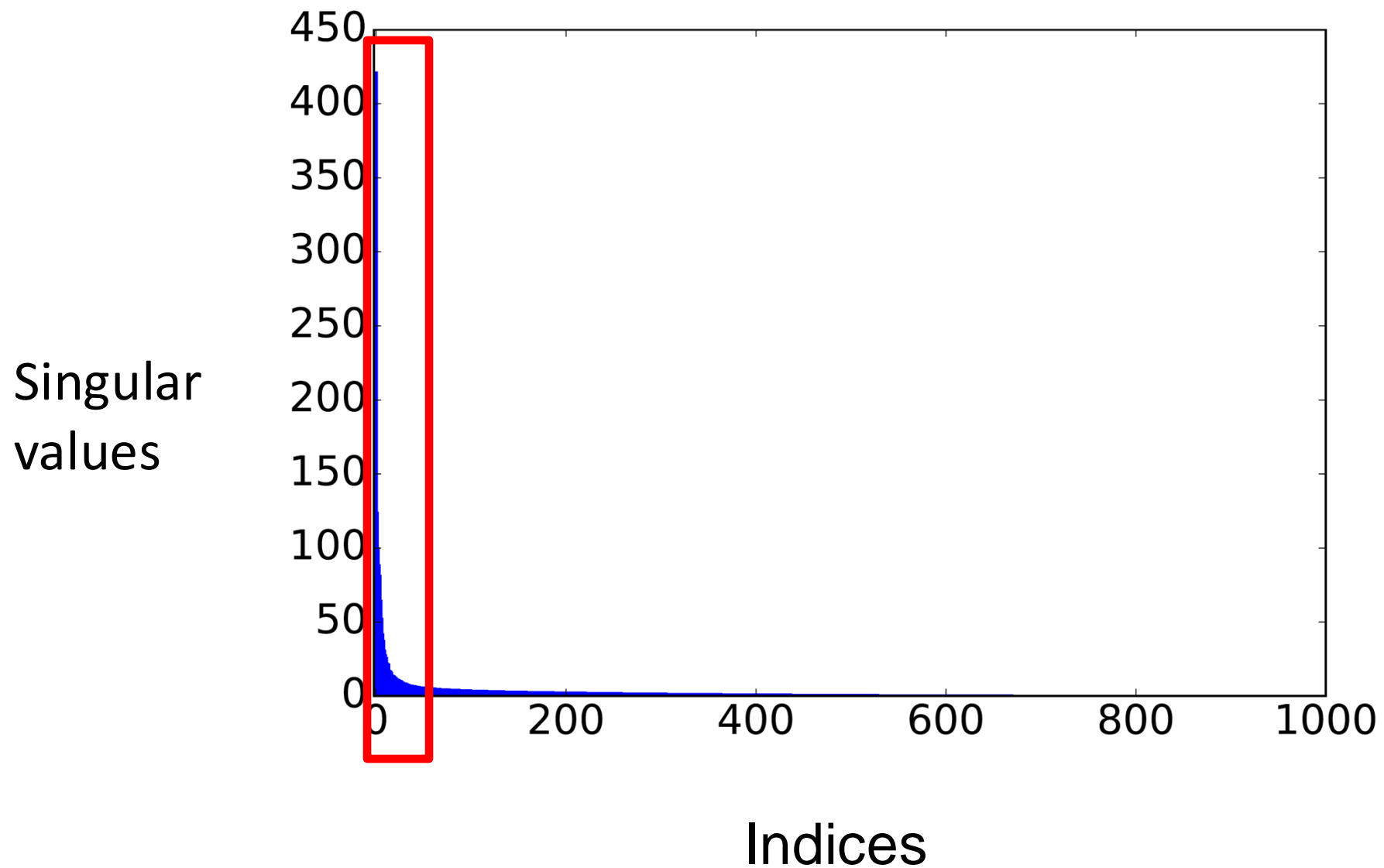


Original image size (1000 x 1500)

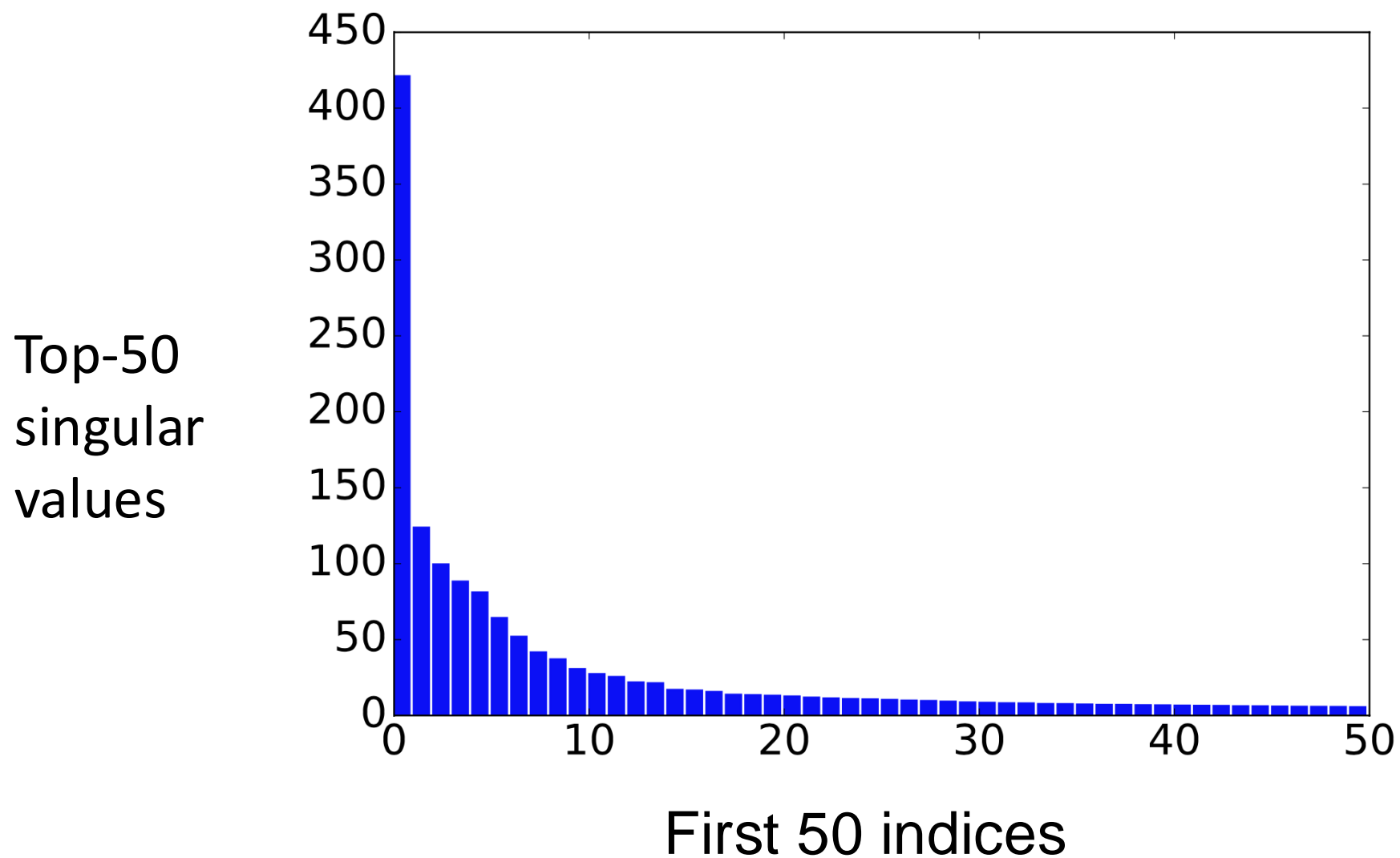
SVD: Example



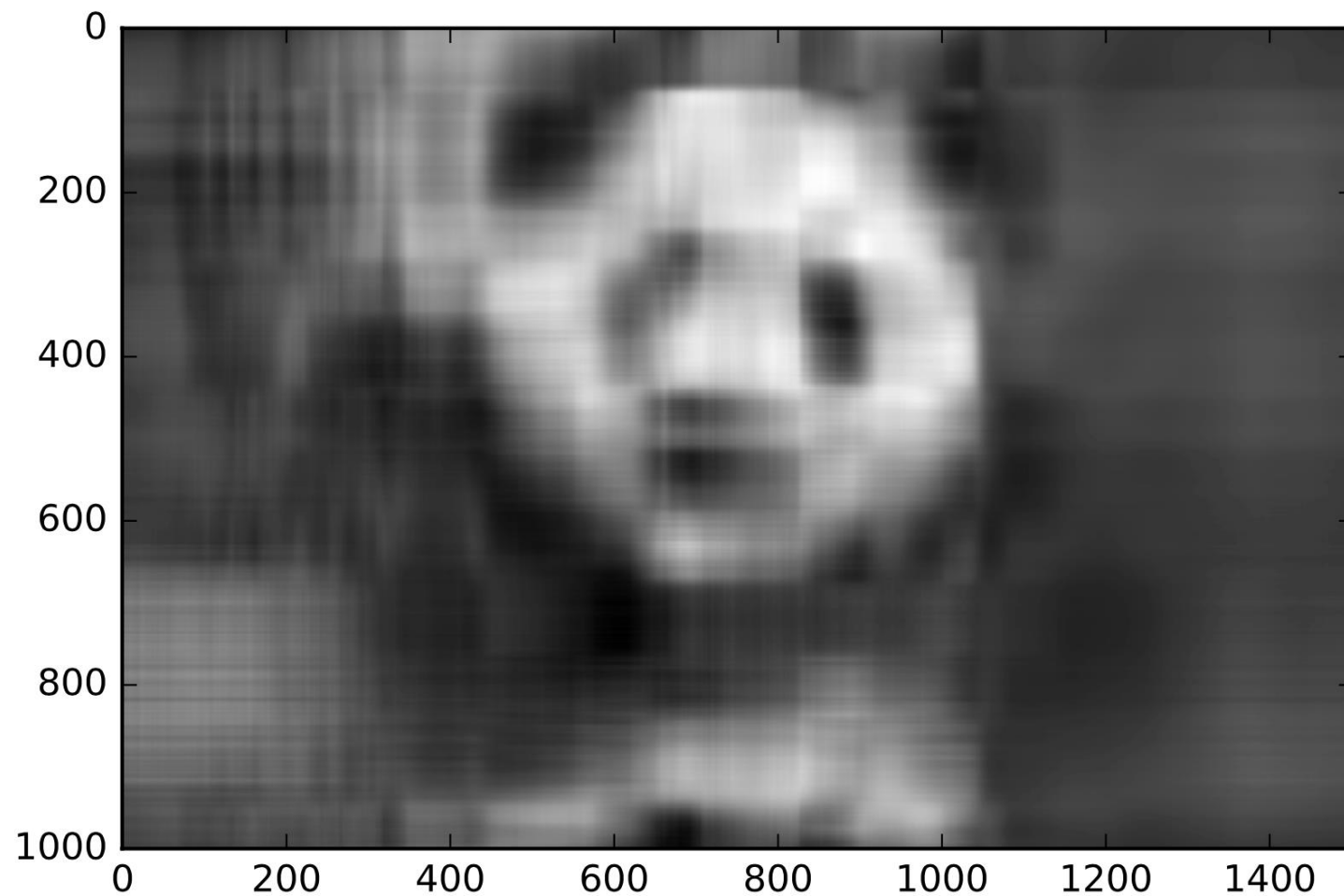
SVD: Example



SVD: Example

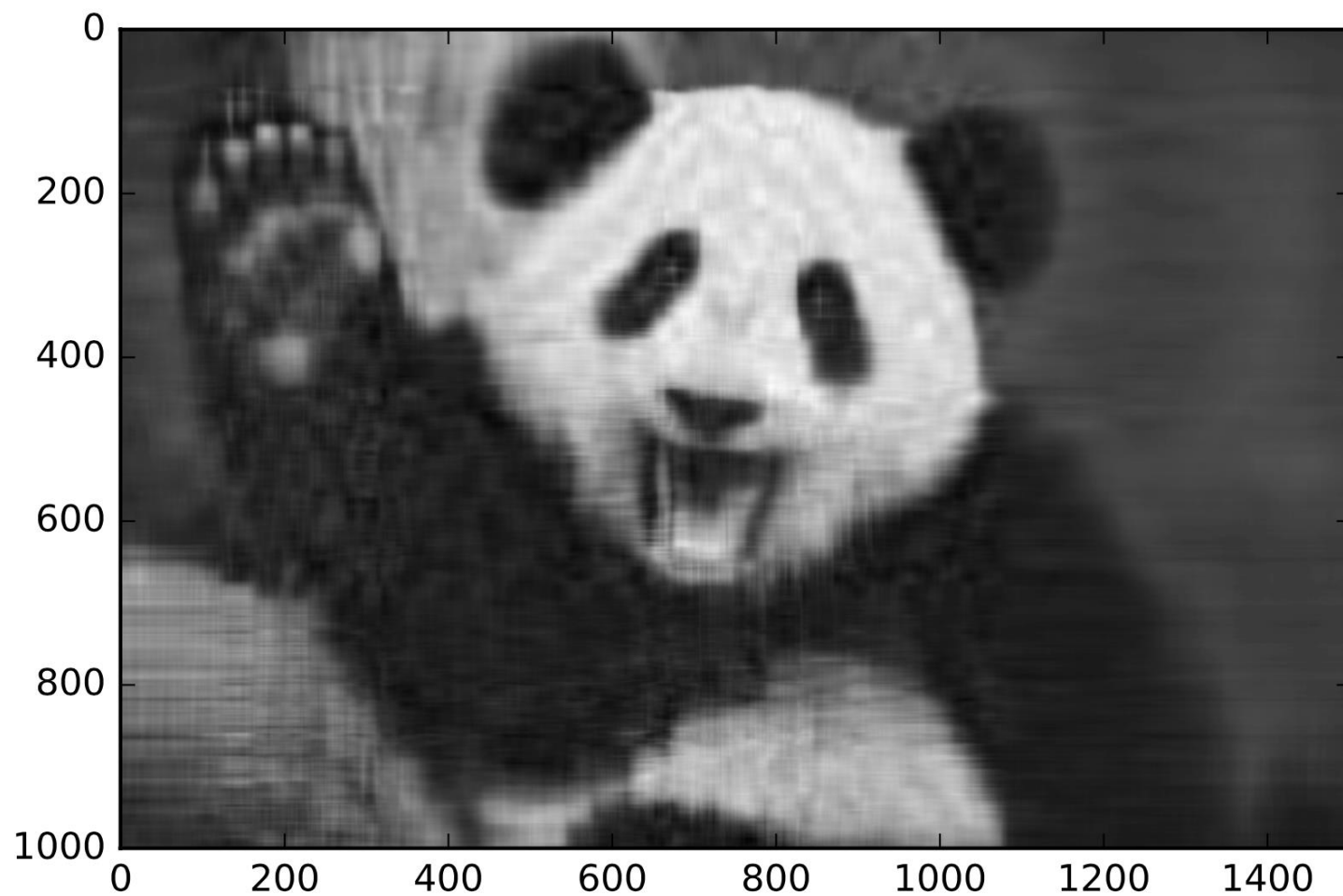


SVD: Example



Rank-5 truncated SVD

SVD: Example



Rank-20 truncated SVD

SVD: Example



Rank-50 truncated SVD

SVD: Example



Rank-100 truncated SVD

SVD: Example



Original image size (1000 x 1500)

SVD: Example



Rank-100 truncated SVD

SVD: Example

- The original matrix
 - Size: 1000×1500
 - #Entries: 1.5M
- The rank-100 truncated SVD
 - $\mathbf{A}_{100} = \sum_{i=1}^{100} \sigma_i \mathbf{u}_i \mathbf{v}_i^T$
 - Size
 - $\{\sigma_i\}$: 100×1
 - $\{\mathbf{u}_i\}$: 100×1000
 - $\{\mathbf{v}_i\}$: 100×1500
 - #Entries: 0.25M
- Truncated SVD saves **83%** storage

Power Iteration for Computing Truncated SVD

A Property

Theorem. If $\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ is the SVD of \mathbf{A} , then $\mathbf{A}^T \mathbf{A} = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$

Proof.

- $\mathbf{A}^T \mathbf{A} = \left(\sum_{i=1}^r \sigma_i \mathbf{v}_i \mathbf{u}_i^T \right) \left(\sum_{j=1}^r \sigma_j \mathbf{u}_j \mathbf{v}_j^T \right)$
- $\mathbf{A}^T \mathbf{A} = \left(\sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{u}_i^T \mathbf{u}_i \mathbf{v}_i^T \right) + \left(\sum_{i \neq j} \sigma_i \sigma_j \mathbf{v}_i \mathbf{u}_i^T \mathbf{u}_j \mathbf{v}_j^T \right)$

$$\left(\sum_i \mathbf{x}_i^T \right) \cdot \left(\sum_j \mathbf{x}_j \right) = \sum_i \mathbf{x}_i^T \mathbf{x}_i + \sum_{i \neq j} \mathbf{x}_i^T \mathbf{x}_j$$

A Property

Theorem. If $A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ is the SVD of A, then $A^T A = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$

Proof.

- $A^T A = \left(\sum_{i=1}^r \sigma_i \mathbf{v}_i \mathbf{u}_i^T \right) \left(\sum_{j=1}^r \sigma_j \mathbf{u}_j \mathbf{v}_j^T \right)$
- $A^T A = \left(\sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{u}_i^T \mathbf{u}_i \mathbf{v}_i^T \right) + \left(\sum_{i \neq j} \sigma_i \sigma_j \mathbf{v}_i \mathbf{u}_i^T \mathbf{u}_j \mathbf{v}_j^T \right)$
- $A^T A = \left(\sum_{i=1}^r \sigma_i^2 \mathbf{v}_i 1 \mathbf{v}_i^T \right) + \left(\sum_{i \neq j} \sigma_i \sigma_j \mathbf{v}_i 0 \mathbf{v}_j^T \right)$

Using the properties of orthonormal basis: $\mathbf{u}_i^T \mathbf{u}_i = 1$ and $\mathbf{u}_i^T \mathbf{u}_j = 0$ for $i \neq j$.

A Property

Theorem. If $\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ is the SVD of \mathbf{A} , then $\mathbf{A}^T \mathbf{A} = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$

Proof.

- $\mathbf{A}^T \mathbf{A} = \left(\sum_{i=1}^r \sigma_i \mathbf{v}_i \mathbf{u}_i^T \right) \left(\sum_{j=1}^r \sigma_j \mathbf{u}_j \mathbf{v}_j^T \right)$
- $\mathbf{A}^T \mathbf{A} = \left(\sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{u}_i^T \mathbf{u}_i \mathbf{v}_i^T \right) + \left(\sum_{i \neq j} \sigma_i \sigma_j \mathbf{v}_i \mathbf{u}_i^T \mathbf{u}_j \mathbf{v}_j^T \right)$
- $\mathbf{A}^T \mathbf{A} = \left(\sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{1} \mathbf{v}_i^T \right) + \left(\sum_{i \neq j} \sigma_i \sigma_j \mathbf{v}_i \mathbf{0} \mathbf{v}_j^T \right)$
- $\mathbf{A}^T \mathbf{A} = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$

A Property

Theorem. If $A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ is the SVD of A , then $A^T A = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$

$$A^T A = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T \quad \longrightarrow \quad A^T A \mathbf{v}_i = \sigma_i^2 \mathbf{v}_i$$

Eigenvalue decomposition
of $A^T A$ to obtain (σ_i, \mathbf{v}_i)

Efficient Power Iteration Method

Power Iteration for Truncated SVD

Goal: Compute the top **1** eigenvalue/eigenvector of $\mathbf{A}^T\mathbf{A} = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$

Algorithm:

1. Randomly initialize a vector \mathbf{x}_0 (with unit ℓ_2 -norm);
2. Repeat the power iteration: $\mathbf{x}_q \leftarrow \mathbf{A}^T\mathbf{A} \mathbf{x}_{q-1}$ and $\mathbf{x}_q \leftarrow \mathbf{x}_q / \|\mathbf{x}_q\|_2$



2 matrix-vector multiplications

$$\mathbf{b} = \mathbf{A} \mathbf{x}_{q-1}$$

$$\mathbf{c} = \mathbf{A}^T \mathbf{b}$$


Cheap computation

Power Iteration for Truncated SVD

Goal: Compute the top **1** eigenvalue/eigenvector of $\mathbf{A}^T \mathbf{A} = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$

Algorithm:

1. Randomly initialize a vector \mathbf{x}_0 (with unit ℓ_2 -norm);
2. Repeat the power iteration: $\mathbf{x}_q \leftarrow \mathbf{A}^T \mathbf{A} \mathbf{x}_{q-1}$ and $\mathbf{x}_q \leftarrow \mathbf{x}_q / \|\mathbf{x}_q\|_2$

Convergence analysis (\mathbf{x}_q converges to \mathbf{v}_1) 

- $\mathbf{x}_0 = \sum_{i=1}^n \alpha_i \mathbf{v}_i$ Every vector can be written as a linear combination of the orthonormal basis
- $\mathbf{x}_q \propto (\mathbf{A}^T \mathbf{A})^q \mathbf{x}_0$ $(\mathbf{A}^T \mathbf{A})^q = \sum_{i=1}^r \sigma_i^{2q} \mathbf{v}_i \mathbf{v}_i^T$
- $\mathbf{x}_q \propto \left(\sum_{i=1}^r \sigma_i^{2q} \mathbf{v}_i \mathbf{v}_i^T \right) \left(\sum_{j=1}^n \alpha_j \mathbf{v}_j \right) = \sum_{i=1}^r \alpha_i \sigma_i^{2q} \mathbf{v}_i$
- $\mathbf{x}_q \propto \sum_{i=1}^r \alpha_i \left(\frac{\sigma_i}{\sigma_1} \right)^{2q} \mathbf{v}_i = \alpha_1 \mathbf{v}_1 + \sum_{i=2}^r \alpha_i \left(\frac{\sigma_i}{\sigma_1} \right)^{2q} \mathbf{v}_i$ Converge to 0 because $\frac{\sigma_i}{\sigma_1} < 1$

Power Iteration for Truncated SVD

Goal: Compute the top k eigenvalue/eigenvector of $\mathbf{A}^T \mathbf{A} = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$

Algorithm:

1. Randomly initialize a vector $\mathbf{X}_0 \in \mathbb{R}^{n \times k}$
 - Entries are i.i.d. standard Gaussian
2. Orthogonalize the columns: $\mathbf{X}_0 \leftarrow \text{orth}(\mathbf{X}_0)$;
3. Repeat the power iteration:
 - i. $\mathbf{X}_q \leftarrow \mathbf{A}^T \mathbf{A} \mathbf{X}_{q-1}$;
 - ii. $\mathbf{X}_q \leftarrow \text{orth}(\mathbf{X}_q)$

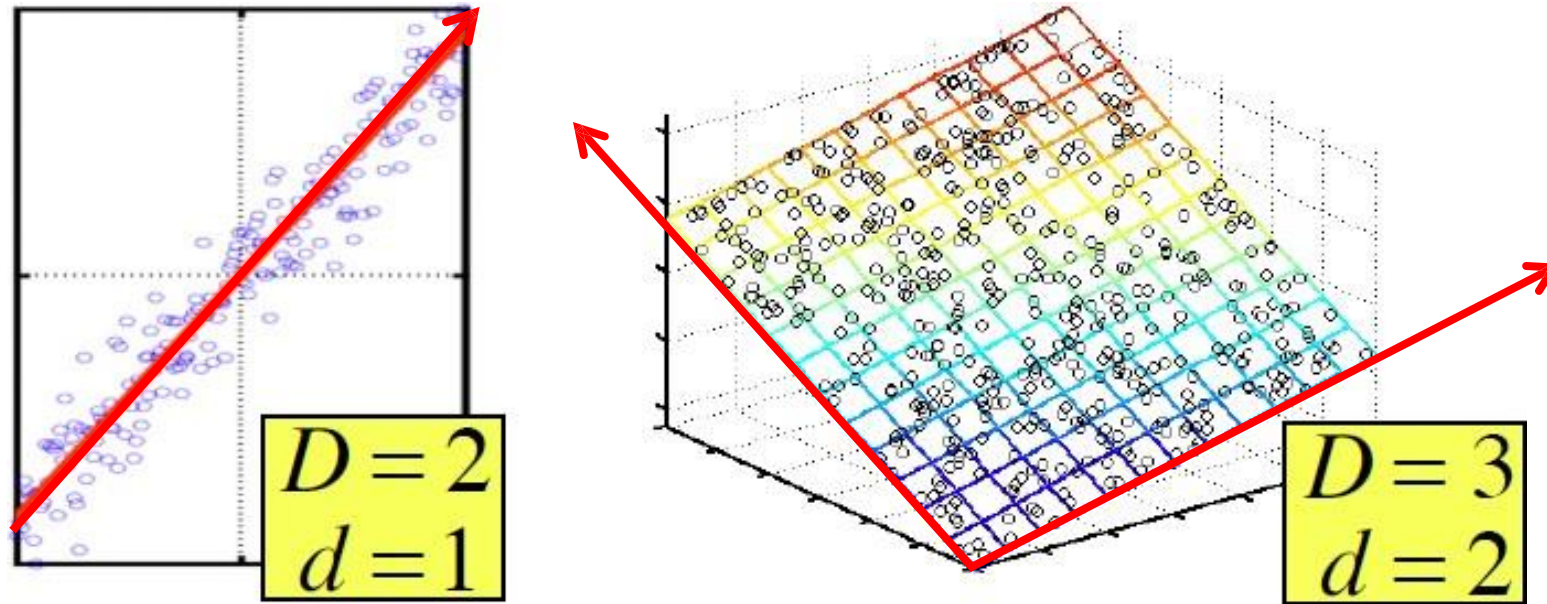
Summary: SVD

- SVD: Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be any matrix
 - $\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$; $r = \text{rank}(\mathbf{A}) \leq \min(m, n)$.
- Truncated SVD: abandon the bottom singular values/vectors
 - $\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$
 - \mathbf{A}_k is the **best** rank- k approximation to \mathbf{A}
- Power iteration (algorithm) for computing truncated SVD

Principal Component Analysis (PCA)

Principal Component Analysis (PCA)

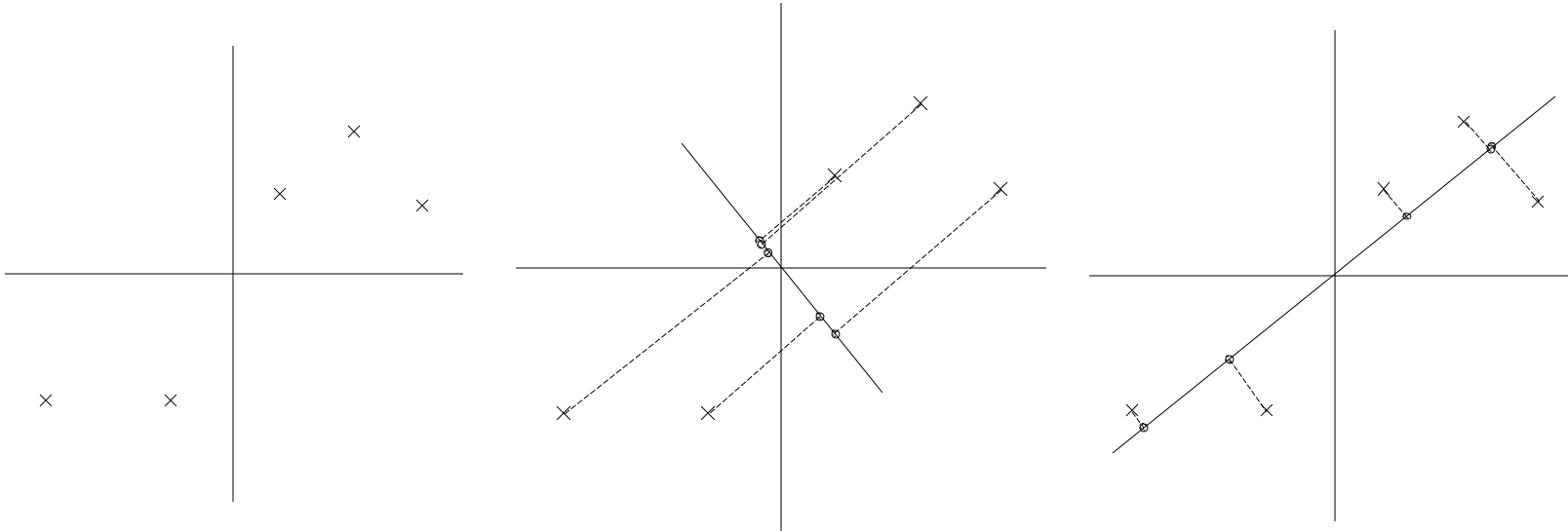
Assumption: Data (approximately) lies on a lower dimensional space



Basis of this subspace are an effective representation of the data

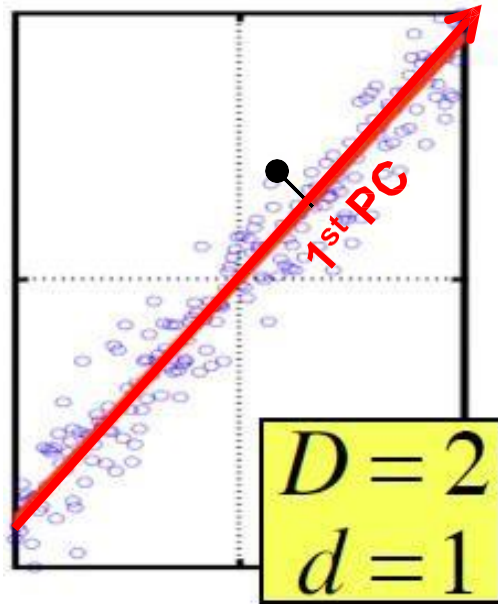
Identifying the basis is known as **Principal Component Analysis**

Which projection is better?



From notes by Andrew Ng

Principal Component Analysis (PCA)



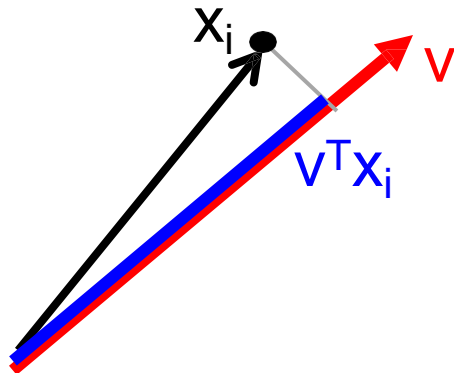
Principal Components (PC) are *orthogonal* directions that capture **most of the variance** in the data

1st PC – direction of greatest variability in data

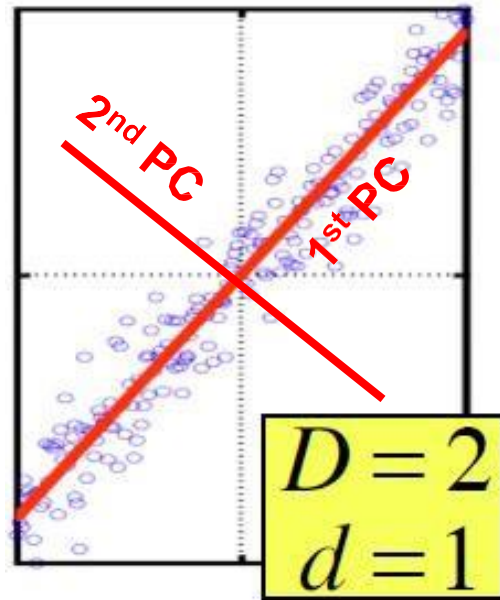
- Projection of data points along 1st PC discriminate the data most along any one direction

Take a data point x_i (D-dimensional vector)

Projection of x_i onto the 1st PC v is $v^T x_i$



Principal Component Analysis (PCA)



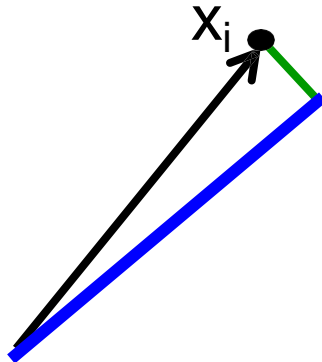
Principal Components (PC) are **orthogonal** directions that capture **most of the variance** in the data

1st PC – direction of greatest variability in data

2nd PC – Next orthogonal direction of greatest variability

- remove all variability in first direction
- then find next direction of greatest variability

And so on ...



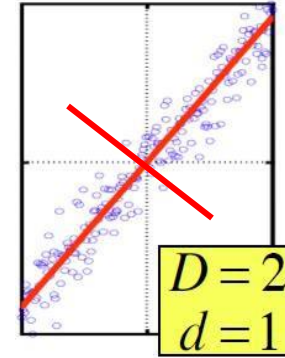
Principal Component Analysis (PCA)

Data points $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$: assume data are centered

Let $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$ denote the principal components

Orthogonal and unit norm $\mathbf{v}_i^T \mathbf{v}_i = 1$ and $\mathbf{v}_i^T \mathbf{v}_j = 0$ for $i \neq j$

Find vector that **maximizes sample variance** of projection



$$\frac{1}{n} \sum_{i=1}^n (\mathbf{v}^T \mathbf{x}_i)^2 = \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v}$$

$$\max_{\mathbf{v}} \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v} \quad \text{s.t.} \quad \mathbf{v}^T \mathbf{v} = 1$$

Lagrangian: $\max_{\mathbf{v}} \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v} - \lambda \mathbf{v}^T \mathbf{v}$

$$\partial / \partial \mathbf{v} = 0$$

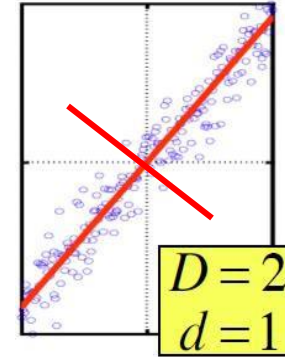
$$(\mathbf{X} \mathbf{X}^T - \lambda \mathbf{I}) \mathbf{v} = 0 \quad \Rightarrow \quad (\mathbf{X} \mathbf{X}^T) \mathbf{v} = \lambda \mathbf{v}$$

Principal Component Analysis (PCA)

$$(XX^T)\mathbf{v} = \lambda\mathbf{v}$$

Therefore, \mathbf{v} is the eigenvector of sample covariance matrix XX^T

$$\text{Sample variance of projection} = \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v} = \lambda \mathbf{v}^T \mathbf{v} = \lambda$$



Thus, eigenvalue λ denotes the amount of variability captured along that eigenvector

Let eigenvalues $\lambda_1 > \lambda_2 > \lambda_3 > \dots$

The **1st PC \mathbf{v}_1** is the **eigenvector** of the sample covariance matrix XX^T associated with the **largest** eigenvalue λ_1

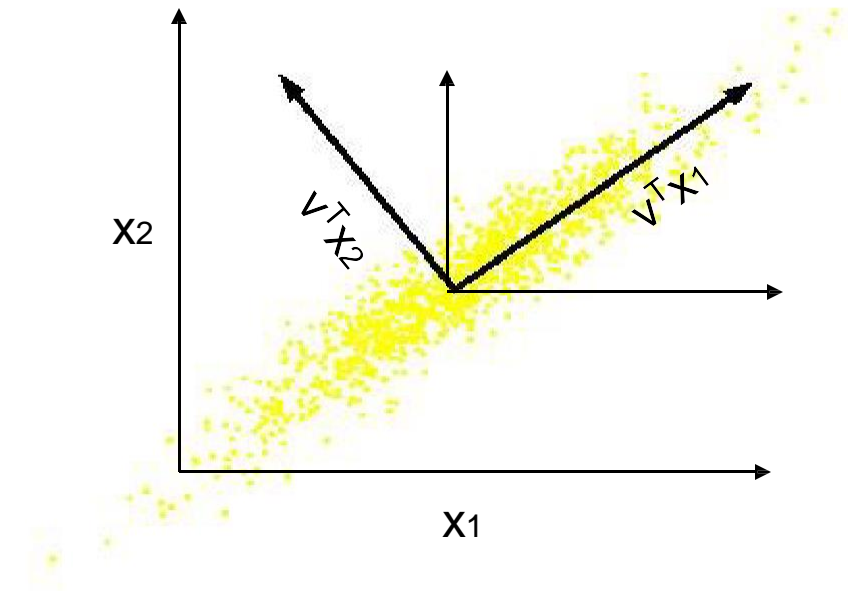
The **2nd PC \mathbf{v}_2** is the **eigenvector** of the sample covariance matrix XX^T associated with the **second largest** eigenvalue λ_2

And so on ...

Computing the Principal Components

The new basis are the eigenvectors of the sample covariance XX^T of the data

Transformed features are uncorrelated



Geometrical interpretation: centering data by **translation** and then followed by **rotation**

- Linear transformation

Another Interpretation

Maximum Variance Subspace

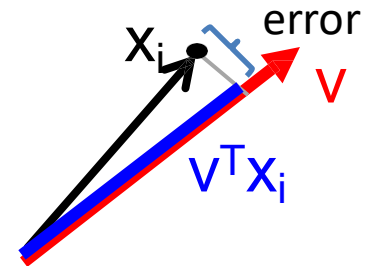
- PCA finds vectors \mathbf{v} such that projections on to the vectors capture **maximum variance in the data**

$$\max_{\mathbf{v}} \frac{1}{n} \sum_{i=1}^n (\mathbf{v}^T \mathbf{x}_i)^2 = \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v}$$

Minimum Reconstruction Error

- PCA finds vectors \mathbf{v} such that projection on to the vectors yields **minimum MSE reconstruction**

$$\min_{\mathbf{v}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i) \mathbf{v}\|^2$$



Dimensionality Reduction using PCA

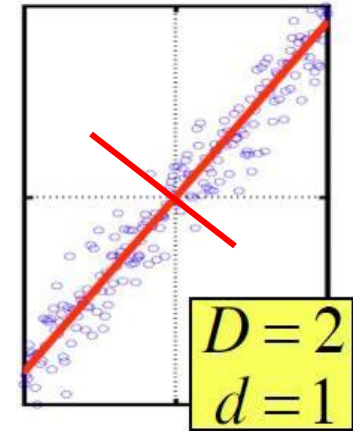
The eigenvalue λ denotes the amount of variability captured along that dimension

Zero eigenvalues indicate no variability along those directions

- Data **exactly** lies on a linear subspace

Keep data projections onto PCs with non-zero eigenvalues

- say v_1, \dots, v_d , where $d = \text{rank}(XX^T)$



Original Representation

Data point in the raw

D -dimensional space

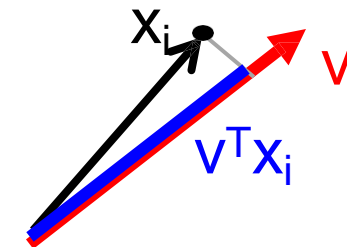
$$\mathbf{x}_i = [x_{i,1}, x_{i,2}, \dots, x_{i,D}]$$

Transformed representation

Projection matrix $\mathbf{V} \in \mathbb{R}^{D \times d}$

(d -dimensional vector)

$$\mathbf{V}^T \mathbf{x}_i = [v_1^T \mathbf{x}_i, v_2^T \mathbf{x}_i, \dots, v_d^T \mathbf{x}_i]$$

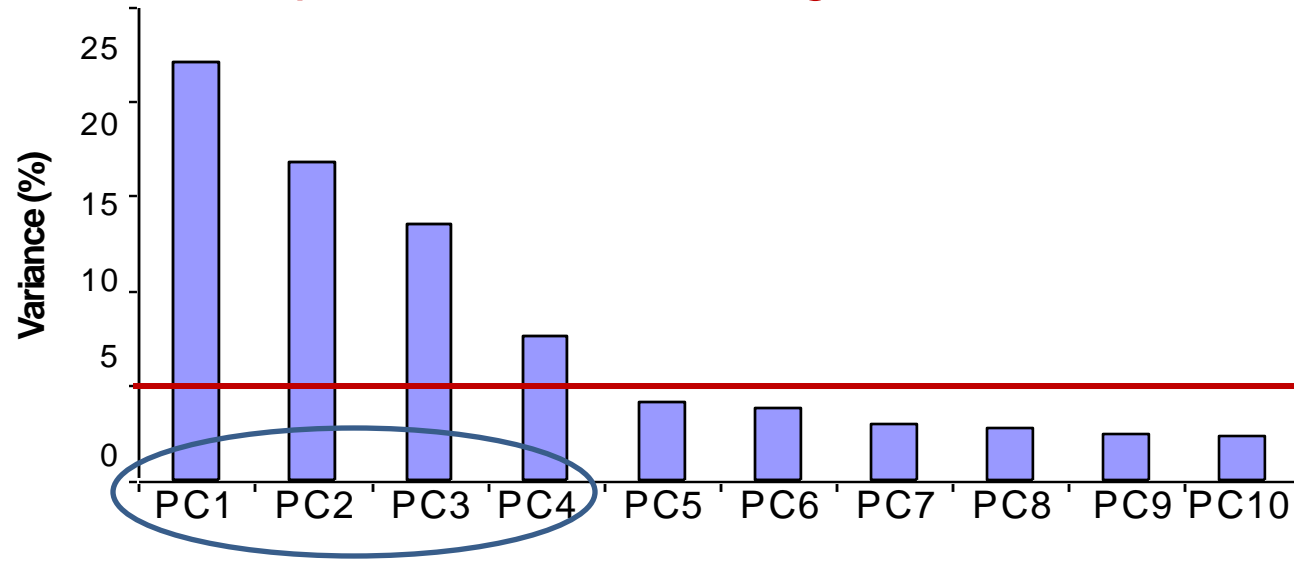


Dimensionality Reduction using PCA

In high-dimensional problem, data usually lies *near* a linear subspace, as noise introduces small variability

Only keep data projections onto PCs with **large** eigenvalues

Can *ignore* the components of lesser significance



Might **lose some information**, but if the eigenvalues are small, don't lose much

Summary: PCA

Project high-dimensional data points into a lower-dimensional space

- **Data points:** $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$
- **Define mean:** $\underline{\mathbf{x}}$, and let $\underline{\mathbf{X}} = \mathbf{X} - \underline{\mathbf{x}}$
- **Principal components:** set of orthonormal basis vectors $(\mathbf{v}_1, \dots, \mathbf{v}_d)$
 - Eigenvalue decomposition: $\underline{\mathbf{X}}\underline{\mathbf{X}}^T \mathbf{v}_i = \lambda_i \mathbf{v}_i$
 - where $\langle \mathbf{v}_j, \mathbf{v}_j \rangle = 1$, and $\langle \mathbf{v}_i, \mathbf{v}_j \rangle = 0$ for $j \neq i$
- **Low-dim representation for \mathbf{x}_i :**
 - $\mathbf{z}_i = (z_{i,1}, \dots, z_{i,d})$, where $z_{i,j} = \langle \mathbf{x}_i - \underline{\mathbf{x}}, \mathbf{v}_j \rangle$

Relationship between SVD and PCA

Centered data points $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{D \times n}$

$$\text{SVD: } \mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^T \quad \mathbf{U}^T\mathbf{U} = \mathbf{I}_r; \mathbf{V}^T\mathbf{V} = \mathbf{I}_r$$

$$\text{PCA: } \mathbf{X}\mathbf{X}^T\mathbf{s}_i = \lambda_i\mathbf{s}_i \quad \longrightarrow \quad \mathbf{X}\mathbf{X}^T\mathbf{S} = \mathbf{S}\Lambda \quad \longrightarrow \quad \mathbf{X}\mathbf{X}^T = \mathbf{S}\Lambda\mathbf{S}^T$$

$$\mathbf{s}_i^T\mathbf{s}_i = 1; \mathbf{s}_i^T\mathbf{s}_j = 0, \forall i \neq j \quad \mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_D] \quad \mathbf{S}^T\mathbf{S} = \mathbf{S}\mathbf{S}^T = \mathbf{I}_D$$

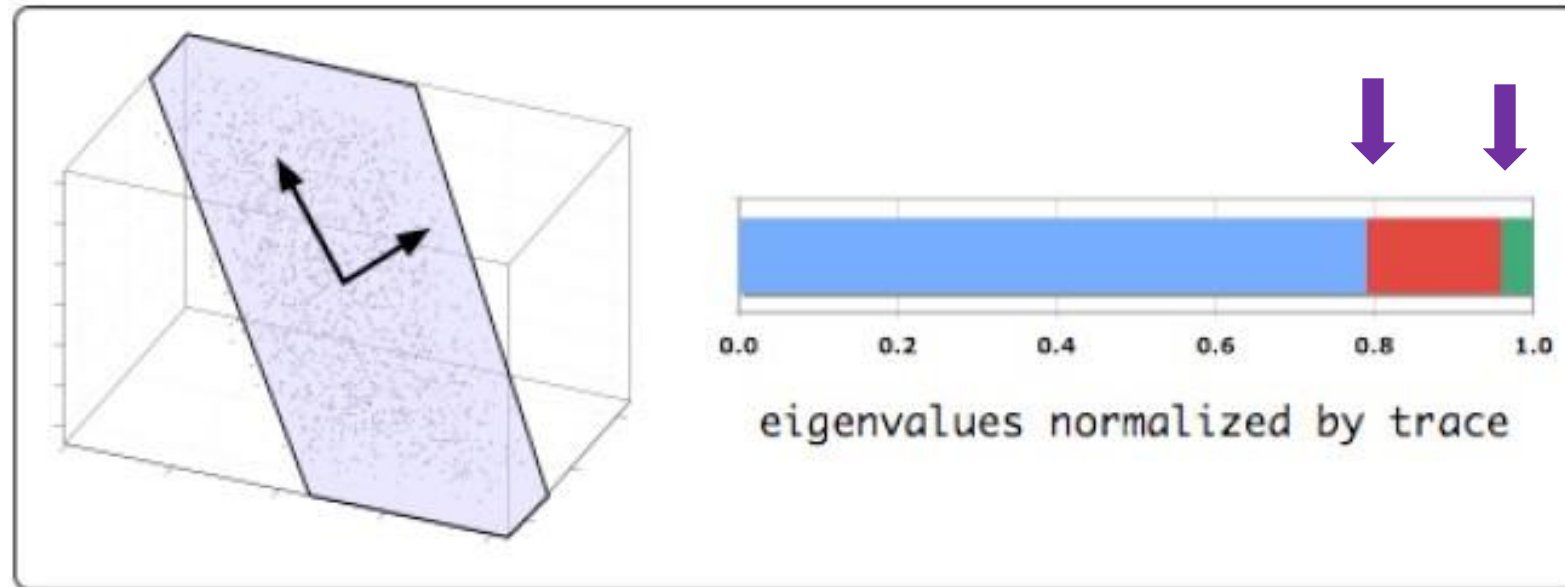
$$\mathbf{X}\mathbf{X}^T = \mathbf{U}\Sigma\mathbf{V}^T\mathbf{V}\Sigma^T\mathbf{U}^T = \mathbf{U}\Sigma^2\mathbf{U}^T$$

The **left singular vectors** (\mathbf{U}) of \mathbf{X} are the same as the **eigenvectors** (\mathbf{S}) of $\mathbf{X}\mathbf{X}^T$

Similarly, the **eigenvalues** (Λ) of $\mathbf{X}\mathbf{X}^T$ are the **squares** of the **singular values** (Σ) of \mathbf{X}

Thus, PCA can reduce to computing the SVD of \mathbf{X} (without forming $\mathbf{X}\mathbf{X}^T$)

Example of PCA



**Eigenvectors and eigenvalues of
covariance matrix for $n=1600$
inputs in $d=3$ dimensions.**

PCA for Face Images

The space of all face images

- Each image as vectors of pixel of values
- Image could be high-dimensional
 - E.g., 100 x 100 image = 10,000 dimensions
- Few 10,000-dim vectors are valid face images
- We want to effectively model the subspace of face images

Eigenfaces [Turk, Pentland '91]

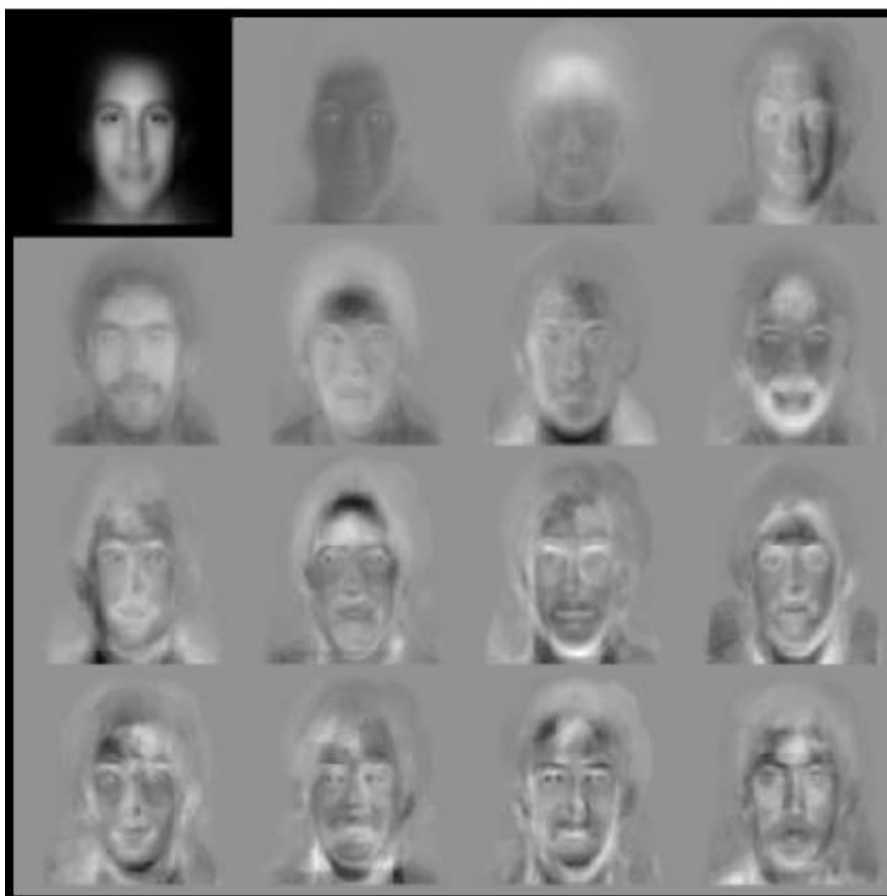
Input images



Eigenfaces:
Principal components



Example: faces



Eigenfaces
from 7562
images:

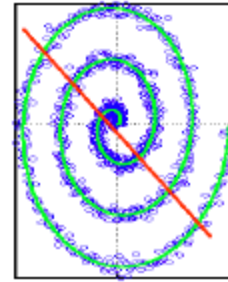
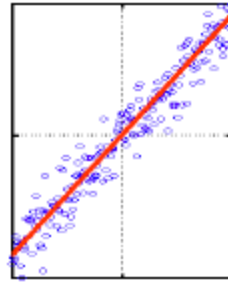
**top left image
is linear
combination
of rest.**

Sirovich & Kirby (1987)
Turk & Pentland (1991)

Properties of PCA

- **Strengths**

- Eigenvector method
- No tuning parameters
- Non-iterative
- No local optima



- **Weaknesses**

- Limited to second order statistics
- Limited to linear projections

Summary

- Singular Value Decomposition (SVD)
 - **Truncated** SVD as dimensionality reduction
 - **Best** rank- k approximation
- Principal Component Analysis (PCA)
 - **Linearly** project high-dim data points into low-dim space
 - Maximize data variance
 - Minimize data mean square reconstruction
- PCA can be reduced to SVD

Acknowledgement

Some slides are from

Arti Singh (CMU)

<https://www.cs.cmu.edu/~aarti/Class/10701/slides/Lecture20.pdf>

Shusen Wang

https://github.com/wangshusen/DeepLearning/blob/master/Slides/5_DR_1.pdf