

Machine Learning for Data Science (CS4786)

Lecture 2

Dimensionality Reduction
&
Principal Component Analysis

Course Webpage :

<http://www.cs.cornell.edu/Courses/cs4786/2016sp/>

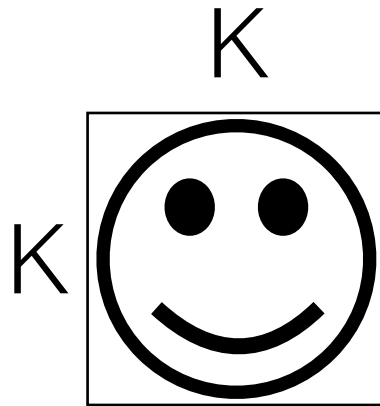
ANNOUNCEMENTS

- Diagnostic assignment due on 4th Feb (Thursday) beginning of class
- Course webpage is the official source of all class related information
- You will be added to CMS once you return Assignment 0 with your net-id on it.

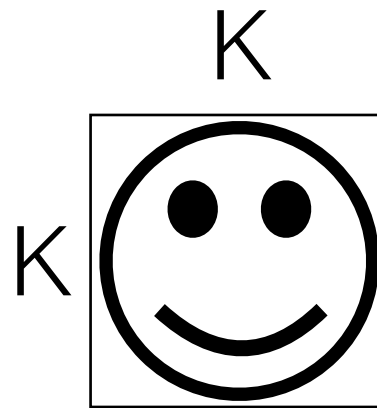
REPRESENTING DATA AS FEATURE VECTORS

- How do we represent data?
- Each data-point often represented as vector referred to as feature vector
- Eg. text document represented by vector in which each coordinate represents a word and value represents number of times the word occurred in the document
- Eg. Image represented as a vector where each coordinate represents a pixel and value represents the grayscale value of that pixel

EXAMPLE: IMAGES

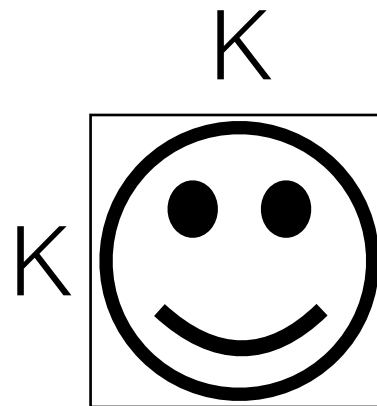


EXAMPLE: IMAGES



vectorize

EXAMPLE: IMAGES



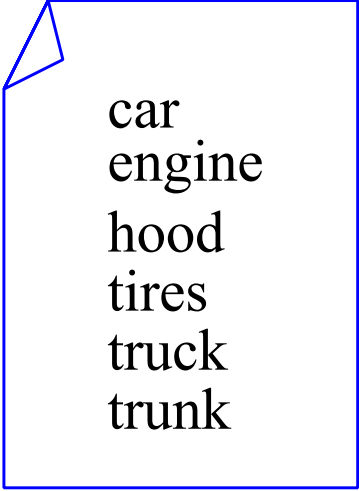
vectorize



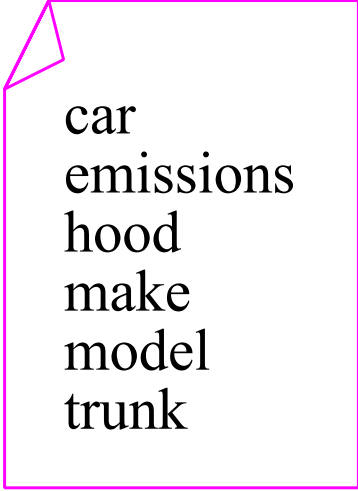
$$d = K^2$$

EXAMPLE: TEXT (BAG OF WORDS)


Documents:



car
engine
hood
tires
truck
trunk



car
emissions
hood
make
model
trunk



Chomsky
corpus
noun
parsing
tagging
wonderful

EXAMPLE: TEXT (BAG OF WORDS)

Documents:

car
engine
hood
tires
truck
trunk

car
emissions
hood
make
model
trunk

Chomsky
corpus
noun
parsing
tagging
wonderful



car	Chomsky	corpus	emissions	engine	hood	make	model	noun	parsing	tagging	tires	truck	trunk	wonderful
1	0	0	0	1	1	0	0	0	0	0	1	1	1	0
1	0	0	1	0	1	1	1	0	0	0	0	0	1	0
0	1	1	0	0	0	0	0	1	1	1	0	0	0	1

DIMENSIONALITY REDUCTION

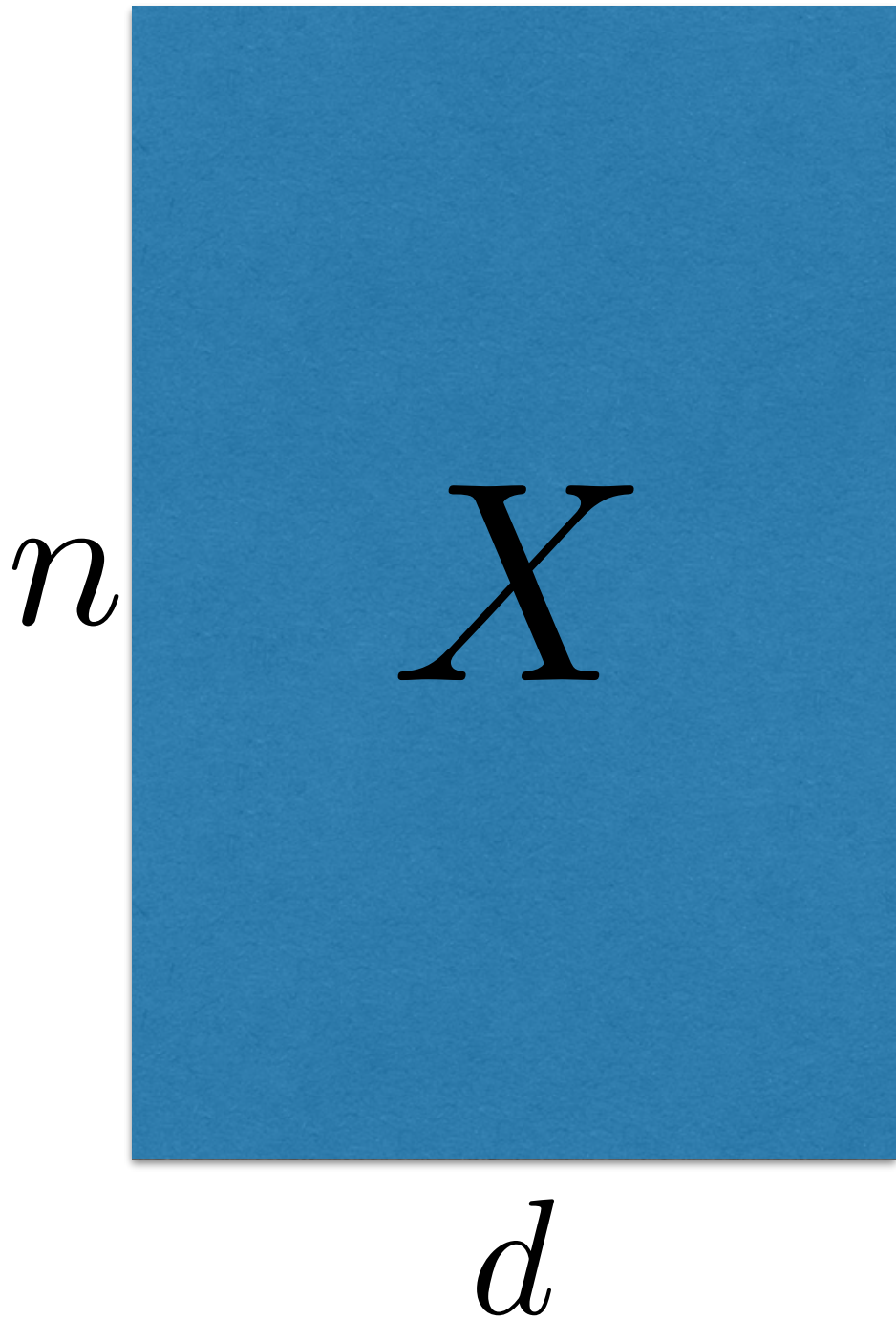
- You are provided with n data points each in \mathbb{R}^d
- Goal: Compress data into n points in \mathbb{R}^K where $K \ll d$
 - Retain as much information about the original data set
 - Retain desired properties of the original data set

DIMENSIONALITY REDUCTION

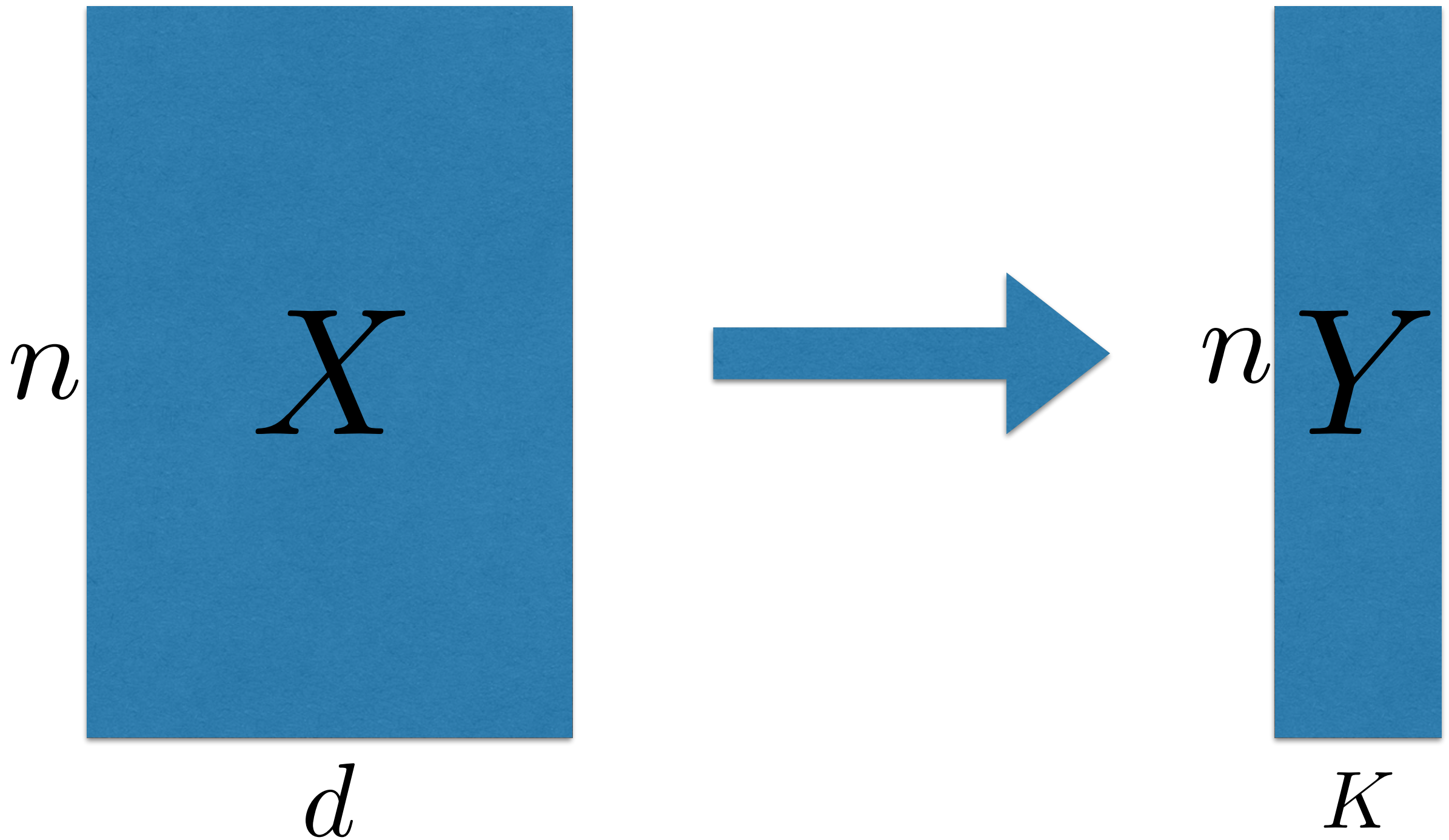
Given feature vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, compress the data points into low dimensional representation $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^K$ where $K \ll d$

DIMENSIONALITY REDUCTION

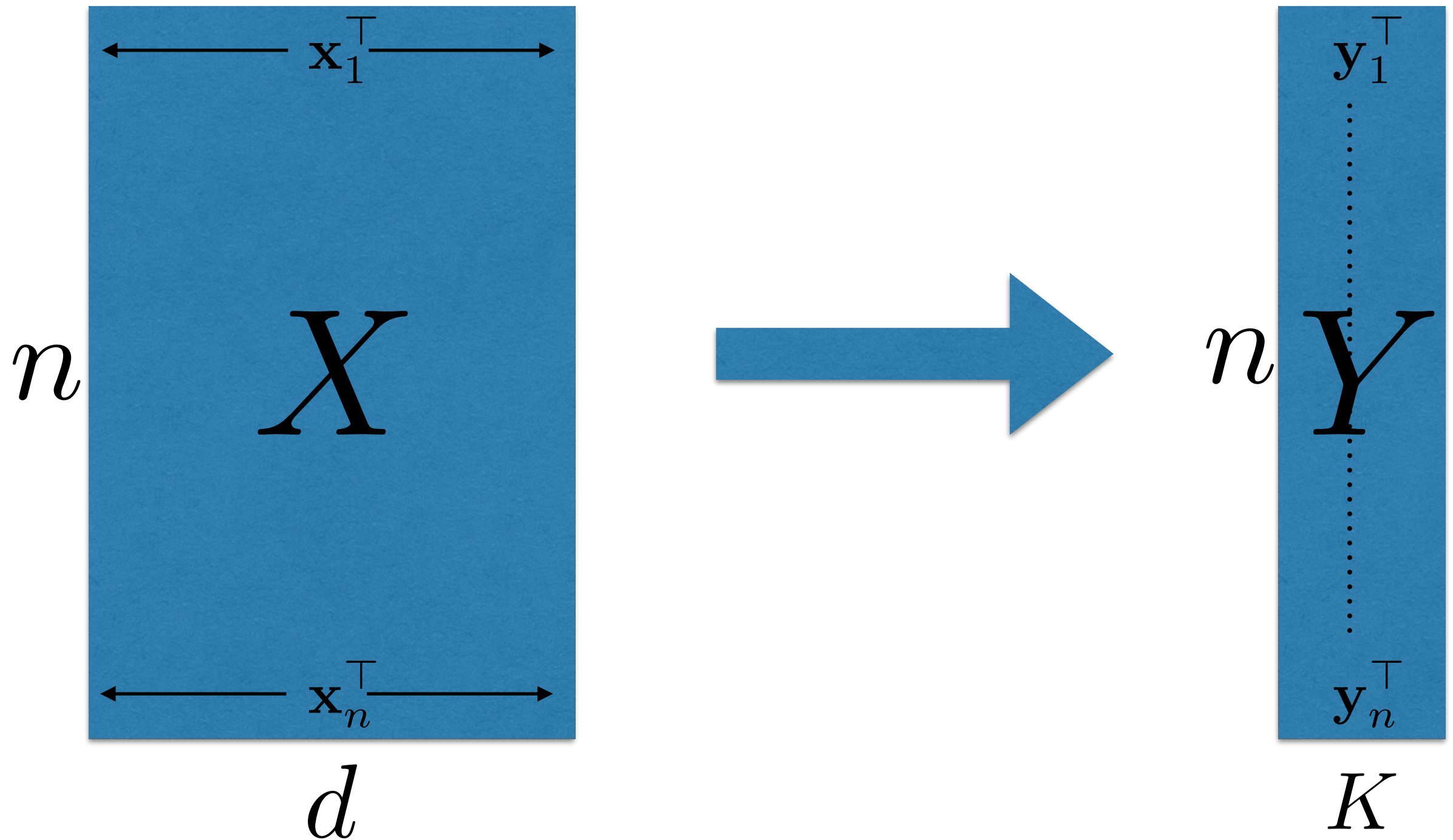
DIMENSIONALITY REDUCTION



DIMENSIONALITY REDUCTION



DIMENSIONALITY REDUCTION



DIMENSIONALITY REDUCTION

Desired properties:

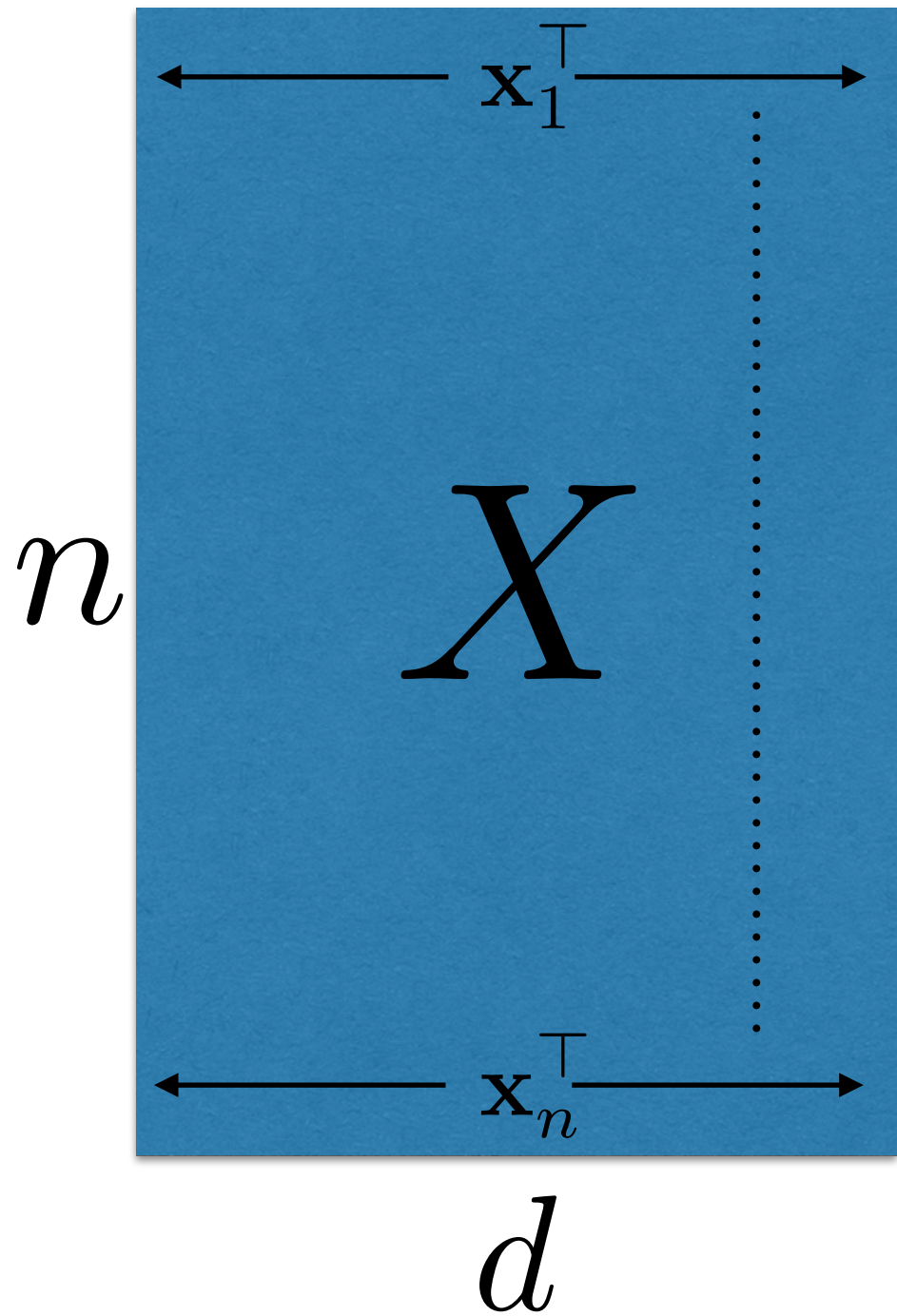
- 1 Original data can be (approximately) reconstructed
- 2 Preserve distances between data points
- 3 “Relevant” information is preserved
- 4 Noise is reduced

DIM REDUCTION: LINEAR TRANSFORMATION

- Pick a low dimensional subspace
- Project linearly to this subspace
- Subspace retains as much information

DIM REDUCTION: LINEAR TRANSFORMATION

DIM REDUCTION: LINEAR TRANSFORMATION



DIM REDUCTION: LINEAR TRANSFORMATION

The diagram illustrates the matrix multiplication $X \times W =$. Matrix X is a blue rectangle with height n and width d . It contains n row vectors, with the first row labeled \mathbf{x}_1^\top and the last row labeled \mathbf{x}_n^\top . A vertical dotted line is on the right side of X . Matrix W is a red rectangle with width d and height K . The equation is shown as $X \times d W =$.

$$\begin{matrix} n \\ \leftarrow \mathbf{x}_1^\top \end{matrix} \begin{matrix} \text{---} X \text{---} \\ \vdots \\ \leftarrow \mathbf{x}_n^\top \end{matrix} \begin{matrix} \text{---} \\ \vdots \\ \text{---} \end{matrix} \begin{matrix} d \\ \text{---} \end{matrix} W \begin{matrix} \text{---} \\ \vdots \\ \text{---} \end{matrix} \begin{matrix} K \\ \text{---} \end{matrix} =$$

DIM REDUCTION: LINEAR TRANSFORMATION

The diagram illustrates a linear transformation for dimensionality reduction. It consists of three main components: a matrix X , a matrix W , and a matrix Y .

- Matrix X :** A large blue rectangle representing an $n \times d$ matrix. The vertical dimension is labeled n on the left. The horizontal dimension is labeled d at the bottom. The top row is labeled \mathbf{x}_1^\top and the bottom row is labeled \mathbf{x}_n^\top . A vertical dotted line is on the right side.
- Matrix W :** A smaller red rectangle representing a $d \times K$ matrix. The horizontal dimension is labeled d on the left. The vertical dimension is labeled K at the bottom.
- Matrix Y :** A blue rectangle representing an $n \times K$ matrix. The vertical dimension is labeled n on the left. The horizontal dimension is labeled K at the bottom. The top row is labeled \mathbf{y}_1^\top and the bottom row is labeled \mathbf{y}_n^\top . A vertical dotted line is on the right side.

The transformation is shown as:

$$X \times W = Y$$

DIM REDUCTION: LINEAR TRANSFORMATION

$$X = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1.1 & 2 & 3 & 4 \\ 3 & 2 & 3 & 4 \\ -1 & 2 & 3 & 4 \\ -0.2 & 2 & 3 & 4 \\ -2 & 2 & 3 & 4 \\ 1.4 & 2 & 3 & 4 \\ 1.4 & 2 & 3 & 4 \\ -0.1 & 2 & 3 & 4 \\ 0.5 & 2 & 3 & 4 \end{bmatrix}$$

DIM REDUCTION: LINEAR TRANSFORMATION

[illegible]

DIM REDUCTION: LINEAR TRANSFORMATION

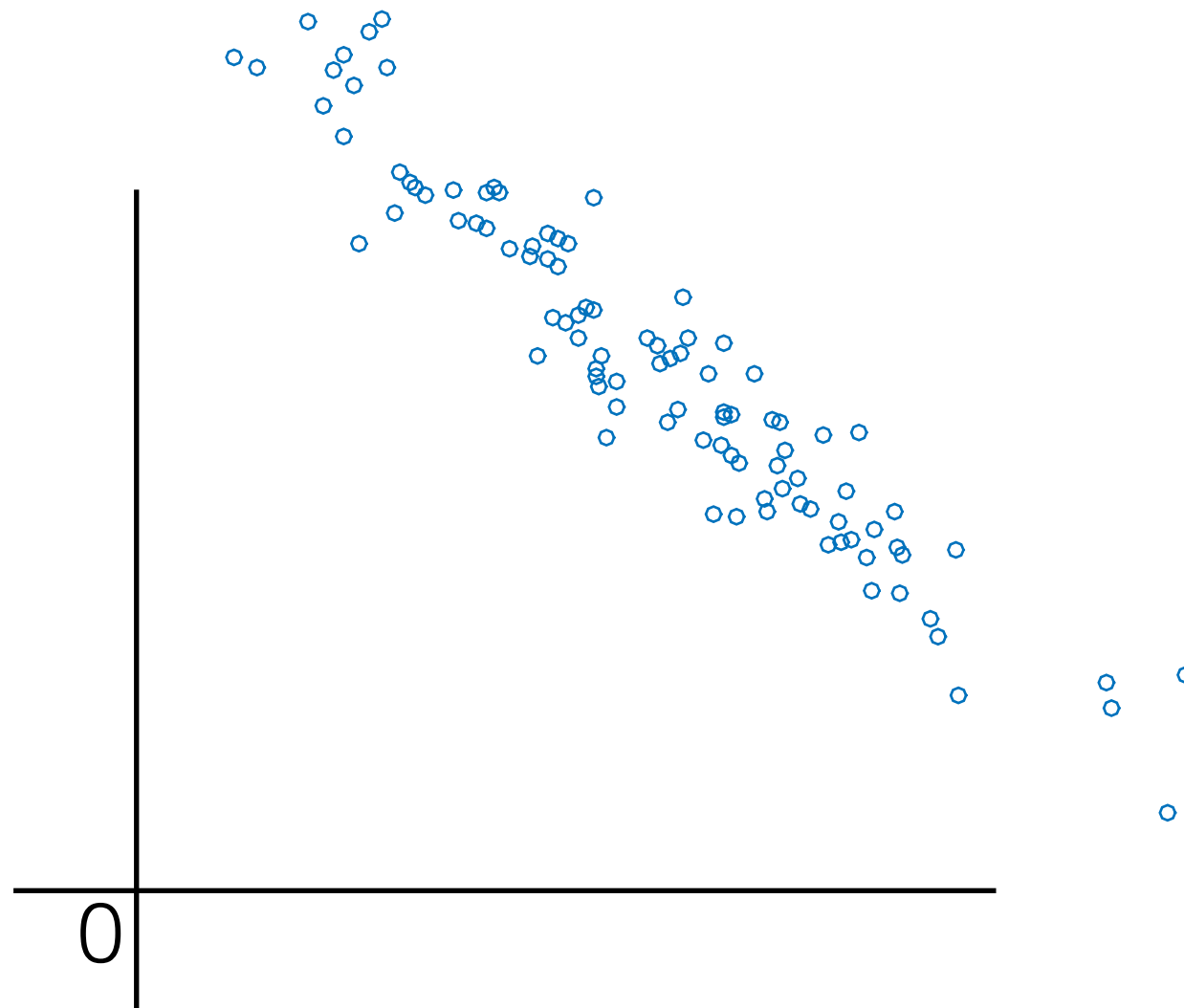
[illegible]

DIM REDUCTION: LINEAR TRANSFORMATION

$$X = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1.1 & 2 & 3 & 4 \\ 3 & 2 & 3 & 4 \\ -1 & 2 & 3 & 4 \\ -0.2 & 2 & 3 & 4 \\ -2 & 2 & 3 & 4 \\ 1.4 & 2 & 3 & 4 \\ 1.4 & 2 & 3 & 4 \\ -0.1 & 2 & 3 & 4 \\ 0.5 & 2 & 3 & 4 \end{bmatrix} = \begin{bmatrix} 1 \\ 1.1 \\ 3 \\ -1 \\ -0.2 \\ -2 \\ 1.4 \\ 1.4 \\ -0.1 \\ 0.5 \end{bmatrix} \times \underset{W}{[1, 0, 0, 0]} + \begin{bmatrix} 0 & 2 & 3 & 4 \\ 0 & 2 & 3 & 4 \\ 0 & 2 & 3 & 4 \\ 0 & 2 & 3 & 4 \\ 0 & 2 & 3 & 4 \\ 0 & 2 & 3 & 4 \\ 0 & 2 & 3 & 4 \\ 0 & 2 & 3 & 4 \\ 0 & 2 & 3 & 4 \\ 0 & 2 & 3 & 4 \end{bmatrix}$$

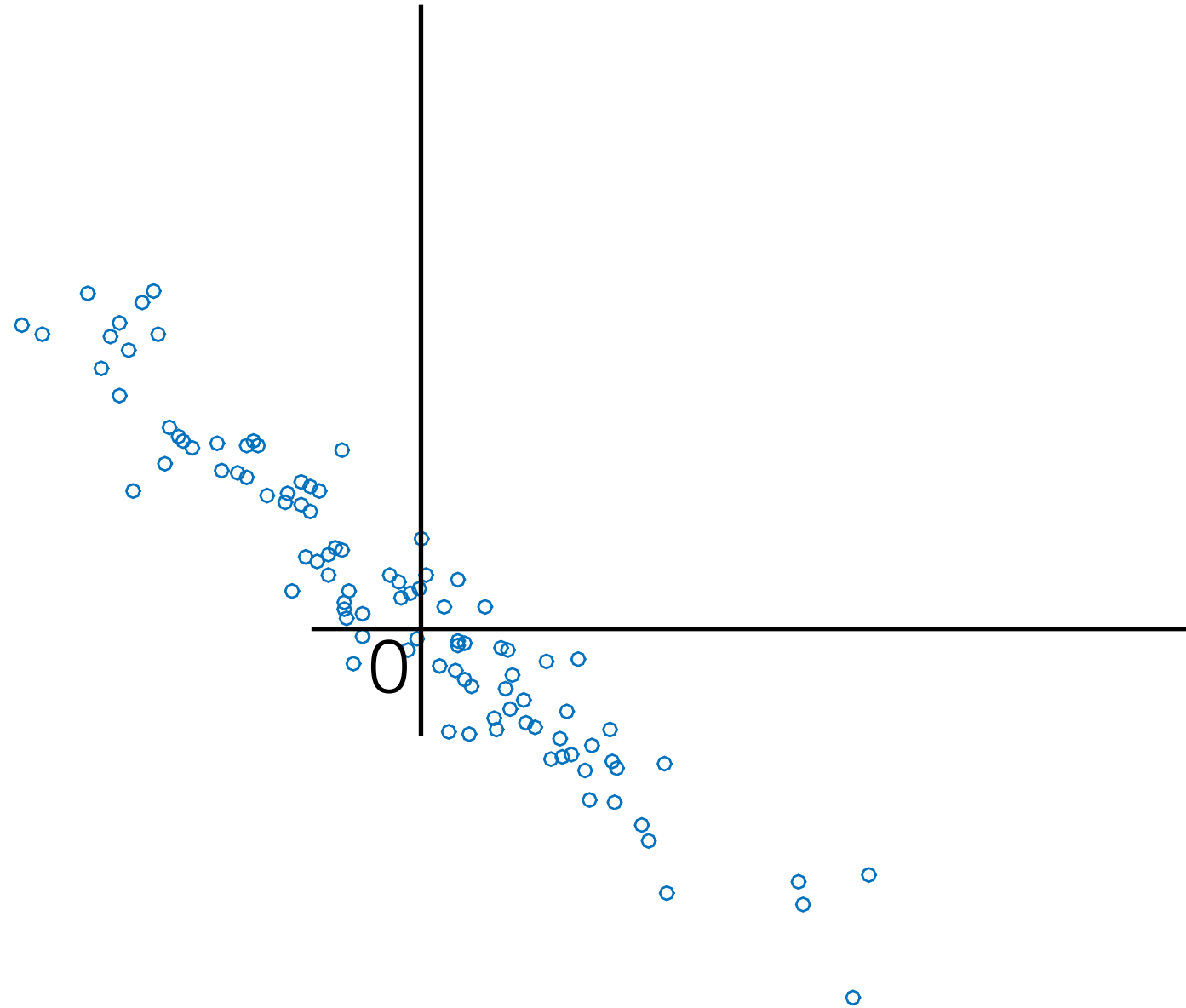
Y

CENTERING DATA



Compressing these data points...

CENTERING DATA



... is same as compressing these.

CENTERING DATA

$$-\mu$$

CENTERING DATA

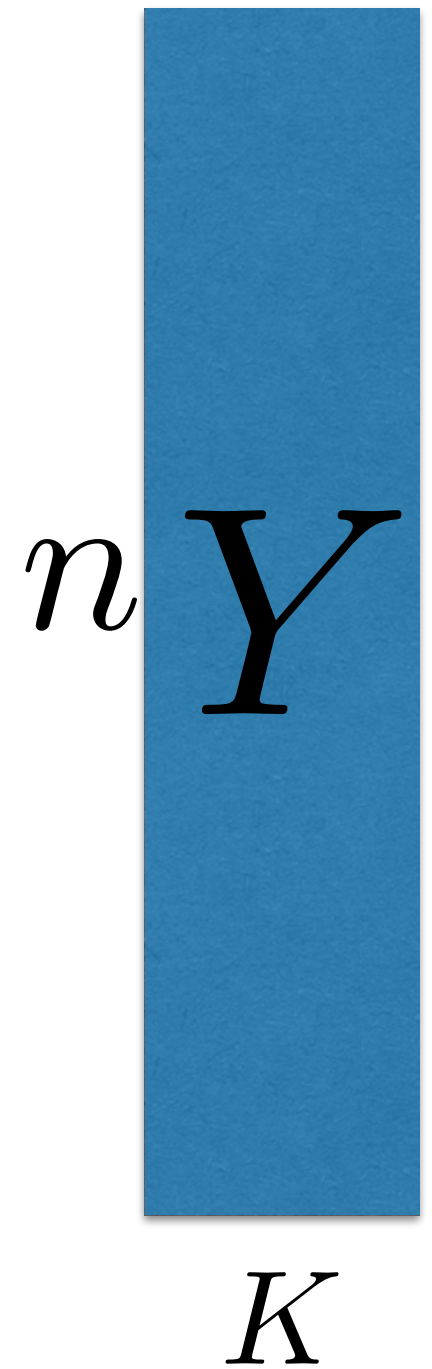
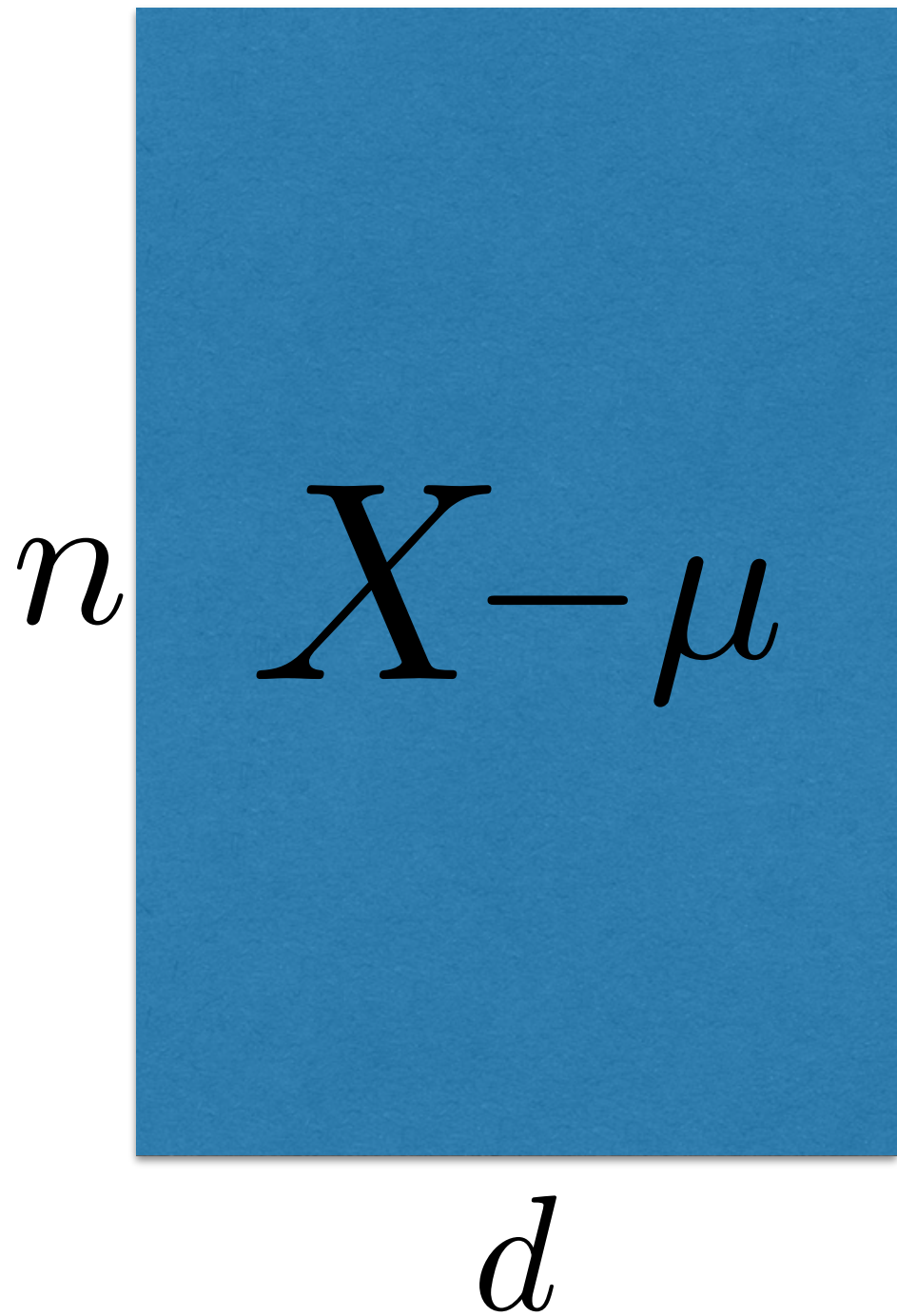


A blue square representing a data matrix. The vertical axis is labeled n and the horizontal axis is labeled d . The matrix is labeled $X - \mu$.

$$n \quad X - \mu$$

d

CENTERING DATA

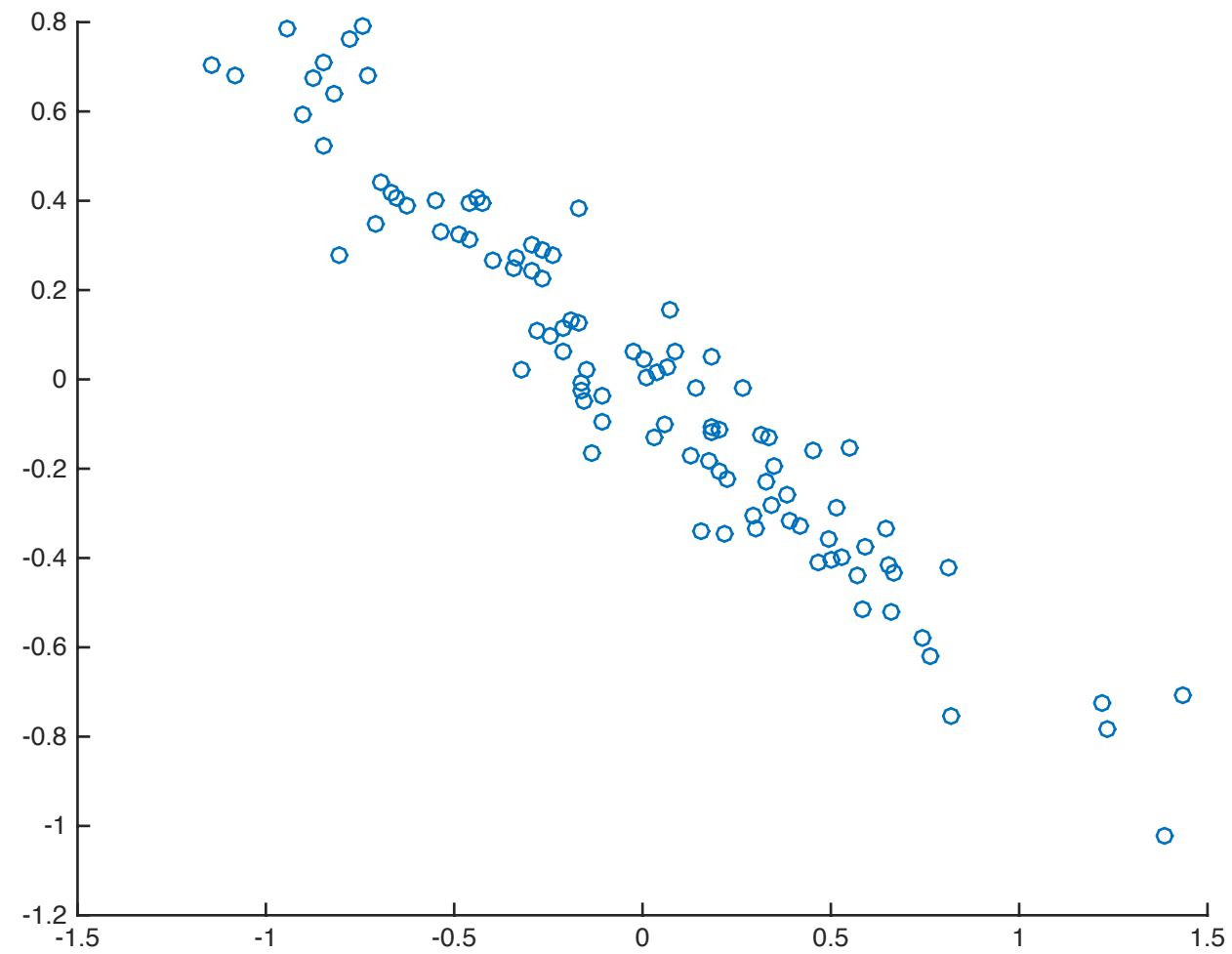


CENTERING DATA

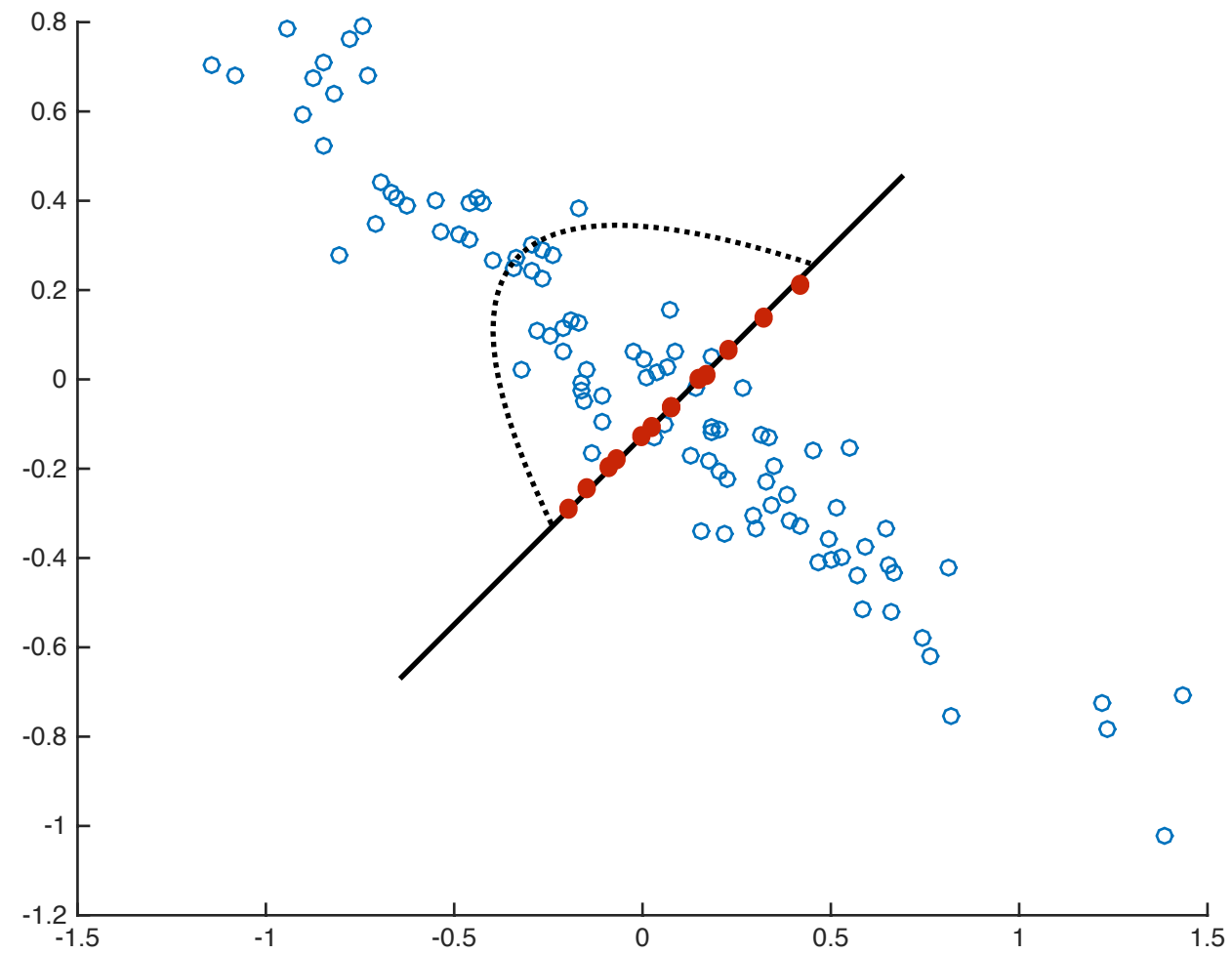
The diagram illustrates the centering of data matrix X . It shows the equation $X - \mu = WY$ using blue rectangular blocks to represent matrices. The first block, representing X , is a square with dimensions n (rows) and d (columns). The second block, representing μ , is a smaller square with dimensions d (rows) and K (columns). The third block, representing W , is a tall, narrow rectangle with dimensions n (rows) and K (columns). The fourth block, representing Y , is a tall, narrow rectangle with dimensions n (rows) and K (columns). The equation is shown as $X - \mu \times W = Y$, with the multiplication symbol \times and the equals sign $=$ placed between the blocks. The dimensions n , d , and K are labeled next to their respective blocks.

$$\begin{matrix} n \\ \times \\ d \end{matrix} X - \mu \times \begin{matrix} d \\ \times \\ K \end{matrix} W = \begin{matrix} n \\ \times \\ K \end{matrix} Y$$

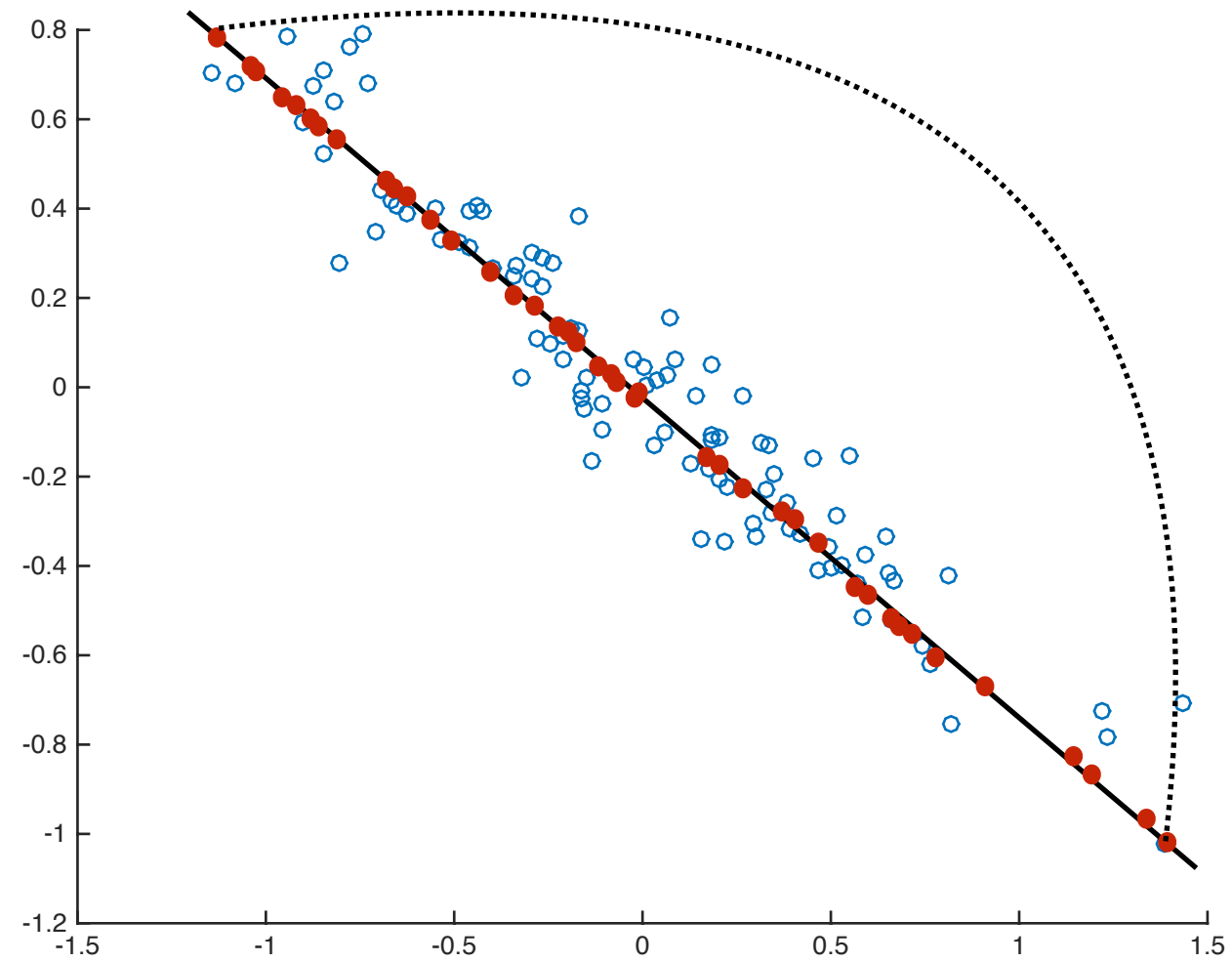
PCA: VARIANCE MAXIMIZATION



PCA: VARIANCE MAXIMIZATION



PCA: VARIANCE MAXIMIZATION



PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most

PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\mathbf{w}_1 = \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2$$

PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\begin{aligned}\mathbf{w}_1 &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top (\mathbf{x}_t - \mu) \right)^2\end{aligned}$$

PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\begin{aligned}\mathbf{w}_1 &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top (\mathbf{x}_t - \mu) \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top (\mathbf{x}_t - \mu) (\mathbf{x}_t - \mu)^\top \mathbf{w}\end{aligned}$$

PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\begin{aligned}\mathbf{w}_1 &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top (\mathbf{x}_t - \mu) \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top (\mathbf{x}_t - \mu) (\mathbf{x}_t - \mu)^\top \mathbf{w} \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \mathbf{w}^\top \Sigma \mathbf{w}\end{aligned}$$

Σ is the covariance matrix

PCA: VARIANCE MAXIMIZATION

Covariance matrix:

$$\Sigma = \frac{1}{n} \sum_{t=1}^n (\mathbf{x}_t - \mu)(\mathbf{x}_t - \mu)^\top$$

PCA: VARIANCE MAXIMIZATION

Covariance matrix:

$$\Sigma = \frac{1}{n} \sum_{t=1}^n (\mathbf{x}_t - \mu)(\mathbf{x}_t - \mu)^\top$$

- Its a $d \times d$ matrix, $\Sigma[i, j]$ measures “covariance” of features i and j

PCA: VARIANCE MAXIMIZATION

Covariance matrix:

$$\Sigma = \frac{1}{n} \sum_{t=1}^n (\mathbf{x}_t - \mu)(\mathbf{x}_t - \mu)^\top$$

- Its a $d \times d$ matrix, $\Sigma[i, j]$ measures “covariance” of features i and j
- Recall $\text{cov}(A, B) = \mathbb{E}[(A - \mathbb{E}[A])(B - \mathbb{E}[B])]$

PCA: VARIANCE MAXIMIZATION

Covariance matrix:

$$\Sigma = \frac{1}{n} \sum_{t=1}^n (\mathbf{x}_t - \mu)(\mathbf{x}_t - \mu)^\top$$

- Its a $d \times d$ matrix, $\Sigma[i, j]$ measures “covariance” of features i and j
- Recall $\text{cov}(A, B) = \mathbb{E}[(A - \mathbb{E}[A])(B - \mathbb{E}[B])]$
- Alternatively,

$$\Sigma[i, j] = \frac{1}{n} \begin{bmatrix} \mathbf{x}_1[j] - \mu[j] \\ \mathbf{x}_2[j] - \mu[j] \\ \dots \\ \mathbf{x}_n[j] - \mu[j] \end{bmatrix}^\top \begin{bmatrix} \mathbf{x}_1[j] - \mu[j] \\ \mathbf{x}_2[j] - \mu[j] \\ \dots \\ \mathbf{x}_n[j] - \mu[j] \end{bmatrix}$$

Inner products measure similarity.

PCA: VARIANCE MAXIMIZATION

- First principal component:

$$\mathbf{w}_1 = \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \mathbf{w}^\top \Sigma \mathbf{w} \quad (1)$$

PCA: VARIANCE MAXIMIZATION

- First principal component:

$$\mathbf{w}_1 = \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \mathbf{w}^\top \Sigma \mathbf{w} \quad (1)$$

To solve the above maximization problem, we use Lagrange multipliers. Specifically there exists λ such that solution \mathbf{w}_1 is:

$$\mathbf{w}_1 = \arg \max_{\mathbf{w}} \mathbf{w}^\top \Sigma \mathbf{w} - \lambda \|\mathbf{w}\|_2^2$$

PCA: VARIANCE MAXIMIZATION

- First principal component:

$$\mathbf{w}_1 = \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \mathbf{w}^\top \Sigma \mathbf{w} \quad (1)$$

To solve the above maximization problem, we use Lagrange multipliers. Specifically there exists λ such that solution \mathbf{w}_1 is:

$$\mathbf{w}_1 = \arg \max_{\mathbf{w}} \mathbf{w}^\top \Sigma \mathbf{w} - \lambda \|\mathbf{w}\|_2^2$$

Taking derivate and equality to 0 we find that $\Sigma \mathbf{w} = \lambda \mathbf{w}$ (ie. eigenvector). Plugging this back into Eq. 1,

$$\frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \Sigma \mathbf{w} = \mathbf{w}^\top (\lambda \mathbf{w}) = \lambda$$

Hence to maximize variance we pick direction with largest eigenvalue

PRINCIPAL COMPONENT ANALYSIS

- Eigenvectors of the covariance matrix are the principal components

PRINCIPAL COMPONENT ANALYSIS

- Eigenvectors of the covariance matrix are the principal components
- Top K principal components are the eigenvectors with K largest eigenvalues

PRINCIPAL COMPONENT ANALYSIS

- Eigenvectors of the covariance matrix are the principal components
- Top K principal components are the eigenvectors with K largest eigenvalues
- $\text{Projection} = \text{Data} \times \text{Top } K \text{ Eigenvectors}$

PRINCIPAL COMPONENT ANALYSIS

- Eigenvectors of the covariance matrix are the principal components
- Top K principal components are the eigenvectors with K largest eigenvalues
- $\text{Projection} = \text{Data} \times \text{Top } K \text{ eigenvectors}$
- $\text{Reconstruction} = \text{Projection} \times \text{Transpose of top } K \text{ eigenvectors}$

PRINCIPAL COMPONENT ANALYSIS

- Eigenvectors of the covariance matrix are the principal components
- Top K principal components are the eigenvectors with K largest eigenvalues
- $\text{Projection} = \text{Data} \times \text{Top } K \text{ eigenvectors}$
- $\text{Reconstruction} = \text{Projection} \times \text{Transpose of top } K \text{ eigenvectors}$
- Independently discovered by Pearson in 1901 and Hotelling in 1933.

PRINCIPAL COMPONENT ANALYSIS: DEMO

