

Dimensionality Reduction

Spring 2024

Hongchang Gao

Why is dimensionality reduction important?

- 1. Exploratory data analysis
 - Visualize high-dimensional data

Car	MPG	Cylinders	Displacement	Horsepower	Weight	Acceleration	Model	Origin
Chevrolet Chevelle	18	8	307	130	3504	12	70	US
Buick Skylark	15	8	350	165	3693	11.5	70	US
Plymouth Satellite	18	8	318	150	3436	11	70	US
AMC Rebel S	16	8	304	150	3433	12	70	US
Ford Torino	17	8	302	140	3449	10.5	70	US
Ford Galaxie	15	8	429	198	4341	10	70	US
Chevrolet Impala	14	8	454	220	4354	9	70	US
Plymouth Fury	14	8	440	215	4312	8.5	70	US
Pontiac Catalina	14	8	455	225	4425	10	70	US
AMC Ambassador	15	8	390	190	3850	8.5	70	US
Citroen DS-2	0	4	133	115	3090	17.5	70	Europe
Chevrolet Chevelle	0	8	350	165	4142	11.5	70	US
Ford Torino (LTD)	0	8	351	153	4034	11	70	US
Plymouth Satellite	0	8	383	175	4166	10.5	70	US
AMC Rebel S	0	8	360	175	3850	11	70	US



Why is dimensionality reduction important?

- 2. Curse of dimensionality
 - Distance concentration refers to the problem of all the pairwise distances between different samples/points in the space converging to the same value as the dimensionality of the data increases.

$$\|x_1 - x_2\|_2 \rightarrow d$$

$$\|x_1 - x_3\|_2 \rightarrow d$$

Why is dimensionality reduction important?

- 3. Computational cost
 - data science algorithms scale linearly with the number of attributes, but very often the scaling is quadratic, or even worse
 - some of the more costly data science algorithms could be impossible to run on larger-sized data.

Why is dimensionality reduction important?

- 4. Noise reduction
 - the real-life data are often noisy and looking at any individual attribute might not provide any insight.
 - if we smartly combine a large number of noisy attributes into a small set of new ones, those new attributes might reveal some useful properties of the data that are not obvious in the original data.

How to reduce the dimensionality?

- How to reduce the dimensionality?
 - Eliminate features?
 - Lose information?

Car	MPG	Cylinders	Displacement	Horsepower	Weight	Acceleration	Model	Origin
Chevrolet Chevelle	18	8	307	130	3504	12	70	US
Buick Skylark	15	8	350	165	3693	11.5	70	US
Plymouth Satellite	18	8	318	150	3436	11	70	US
AMC Rebel S	16	8	304	150	3433	12	70	US
Ford Torino	17	8	302	140	3449	10.5	70	US
Ford Galaxie	15	8	429	198	4341	10	70	US
Chevrolet Impala	14	8	454	220	4354	9	70	US
Plymouth Fury	14	8	440	215	4312	8.5	70	US
Pontiac Catalina	14	8	455	225	4425	10	70	US
AMC Ambassador	15	8	390	190	3850	8.5	70	US
Citroen DS-2	0	4	133	115	3090	17.5	70	Europe
Chevrolet Chevelle	0	8	350	165	4142	11.5	70	US
Ford Torino (L)	0	8	351	153	4034	11	70	US
Plymouth Satellite	0	8	383	175	4166	10.5	70	US
AMC Rebel S	0	8	360	175	3850	11	70	US

Preserve information when reducing the dimensionality!

Toy Example

- 3D points
 - If each component is stored in a byte, 18 bytes are needed

1	2	4	3	5	6
2	4	8	6	10	12
3	6	12	9	15	18

Toy Example

- All points are correlated: a common point scaled by different factors

<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3	= 1 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3	<table><tr><td>2</td></tr><tr><td>4</td></tr><tr><td>6</td></tr></table>	2	4	6	= 2 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3	<table><tr><td>4</td></tr><tr><td>8</td></tr><tr><td>12</td></tr></table>	4	8	12	= 4 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3
1																										
2																										
3																										
1																										
2																										
3																										
2																										
4																										
6																										
1																										
2																										
3																										
4																										
8																										
12																										
1																										
2																										
3																										
<table><tr><td>3</td></tr><tr><td>6</td></tr><tr><td>9</td></tr></table>	3	6	9	= 3 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3	<table><tr><td>5</td></tr><tr><td>10</td></tr><tr><td>15</td></tr></table>	5	10	15	= 5 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3	<table><tr><td>6</td></tr><tr><td>12</td></tr><tr><td>18</td></tr></table>	6	12	18	= 6 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3
3																										
6																										
9																										
1																										
2																										
3																										
5																										
10																										
15																										
1																										
2																										
3																										
6																										
12																										
18																										
1																										
2																										
3																										

Toy Example

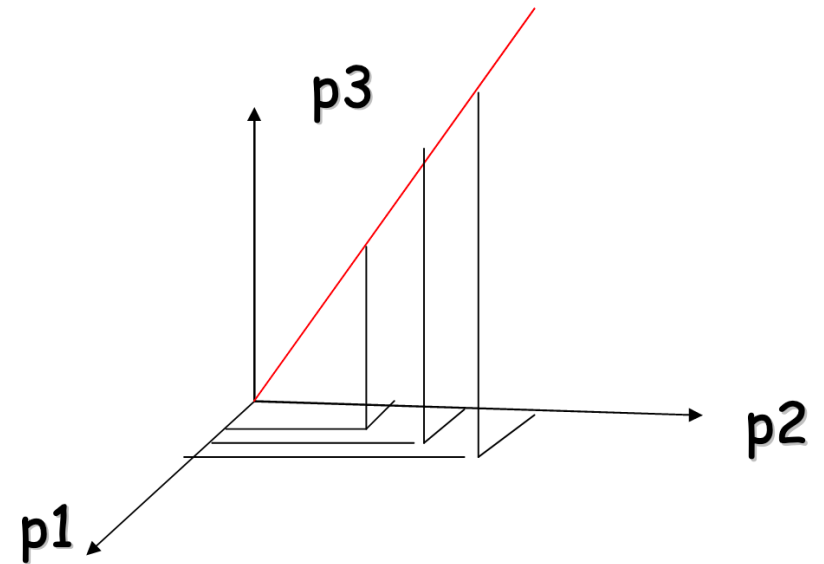
- To store them, only 9 bytes are needed
 - Store one point and six scaling factors

1		1		2		1		4		1
2	= 1 *	2		4	= 2 *	2		8	= 4 *	2
3		3		6		3		12		3
3		1		5		1		6		1
6	= 3 *	2		10	= 5 *	2		12	= 6 *	2
9		3		15		3		18		3

Toy Example

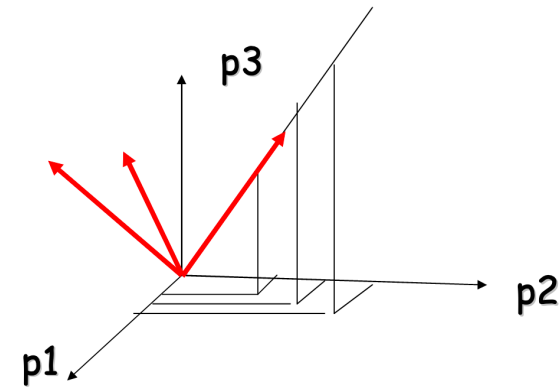
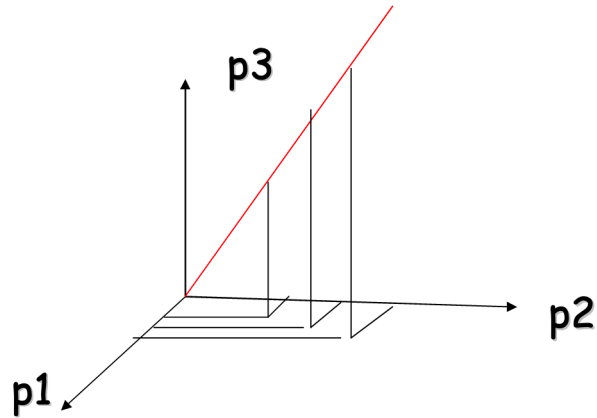
- Geometrical interpretation
 - All data points lie in a 1D subspace of the original 3D space

1	2	4	3	5	6
2	4	8	6	10	12
3	6	12	9	15	18



Toy Example

- Geometrical interpretation
 - Find a new coordinate system where one of the axes is along the direction of the line
 - In the new coordinate system, every data point has only one non-zero coordinate



Principal Components Analysis (PCA)

- PCA is a dimensionality reduction method
 - A linear transformation
 - Find a new coordinate system for the dataset
 - Only use **a small part of coordinates** to represent data points
 - Preserve as much of the data's variance as possible
- Formally, given a dataset with n samples $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ $\mathbf{x}_i \in \mathbb{R}^d$
 - Find a **linear transformation** $W^{d \times k}$ where $k < d$
 - d is the number of features in the original data
 - k is the number of new features
 - Preserve the variance as much as possible

Principal Components Analysis

- First principal component

- Subtract the mean

$$\tilde{X} = [\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_n] \quad \tilde{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$$

- Apply the linear transformation
 - Get data points in the new coordinate system

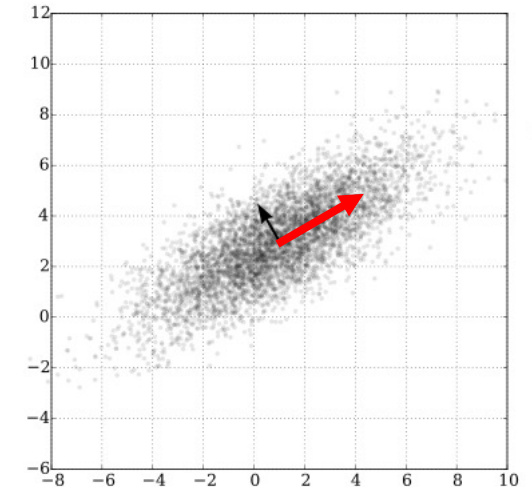
$$[\mathbf{w}^T \tilde{\mathbf{x}}_1, \mathbf{w}^T \tilde{\mathbf{x}}_2, \dots, \mathbf{w}^T \tilde{\mathbf{x}}_n]$$

- Compute the **variance** in the new coordinate system

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

$$\mathbf{x}=(1, 0), \mathbf{y}=(0, 1), \\ \mathbf{a}=(2, 3)$$

$$\mathbf{x}^T \mathbf{a} = 2, \mathbf{y}^T \mathbf{a} = 3$$



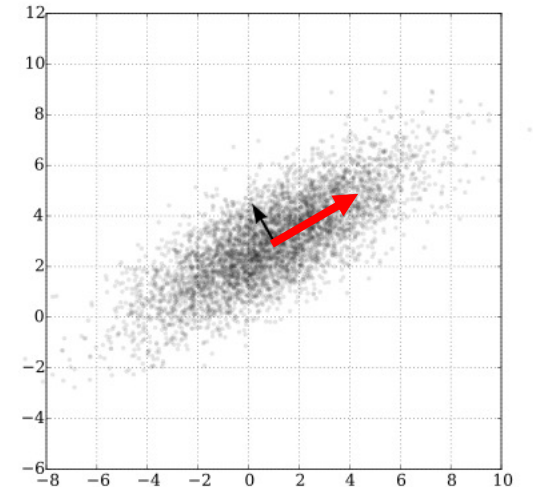
Principal Components Analysis

- How to find this new coordinate system/linear transformation?
 - Maximize the variance in the new coordinate system

$$\max_{\|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \tilde{\mathbf{x}}_i)^2$$

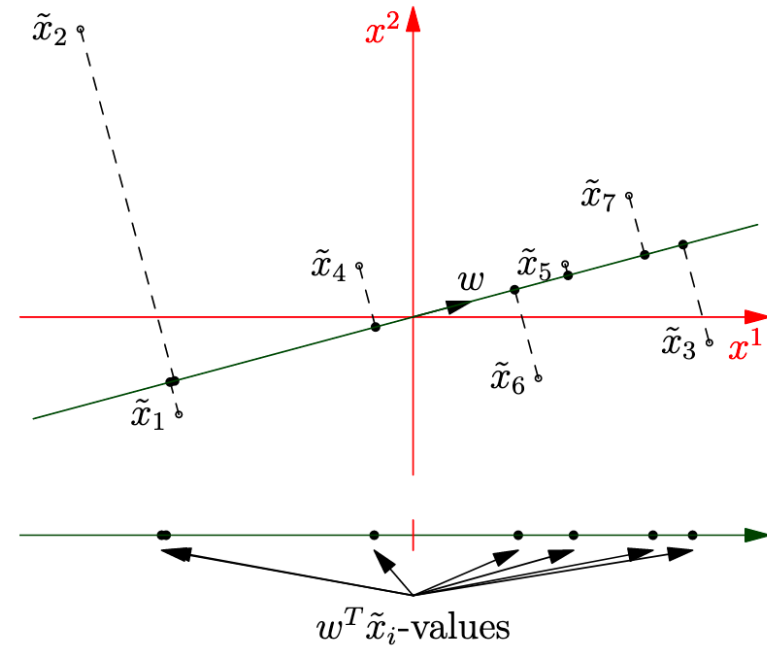
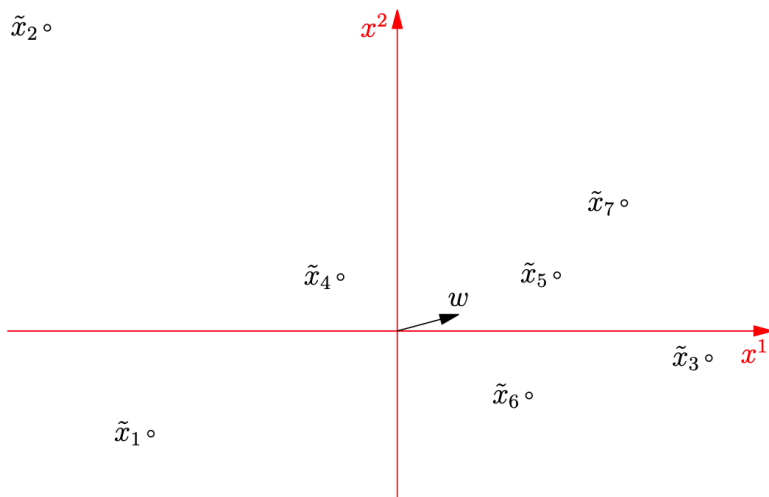
- Optimizing this problem can find the new coordinate system
 - The first principal component: $d \rightarrow 1$

$$\{\mathbf{w}^T \tilde{\mathbf{x}}_1, \mathbf{w}^T \tilde{\mathbf{x}}_2, \dots, \mathbf{w}^T \tilde{\mathbf{x}}_n\}$$



Principal Components Analysis

- Visualization



Principal Components Analysis

- Second principal component

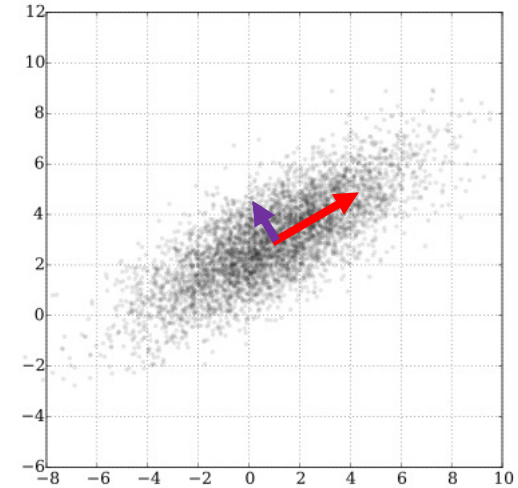
$$\max_{\substack{\|\mathbf{w}\|_2=1, \\ \mathbf{w} \perp \mathbf{w}_1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \tilde{\mathbf{x}}_i)^2$$

- The first principal component

$$\{\mathbf{w}_1^T \tilde{\mathbf{x}}_1, \mathbf{w}_1^T \tilde{\mathbf{x}}_2, \dots, \mathbf{w}_1^T \tilde{\mathbf{x}}_n\}$$

- The second principal component

$$\{\mathbf{w}_2^T \tilde{\mathbf{x}}_1, \mathbf{w}_2^T \tilde{\mathbf{x}}_2, \dots, \mathbf{w}_2^T \tilde{\mathbf{x}}_n\}$$



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

$$W = [\mathbf{w}_1, \mathbf{w}_2] \in \mathbb{R}^{d \times 2}$$



$$\hat{X} = W^T \tilde{X} \in \mathbb{R}^{2 \times n}$$

Principal Components Analysis

- The k-th principal component

$$\max_{\substack{\|\mathbf{w}\|_2=1, \\ \mathbf{w} \perp \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{k-1}}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \tilde{\mathbf{x}}_i)^2$$

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

$$W = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k] \in \mathbb{R}^{d \times k}$$



$$\hat{X} = W^T \tilde{X} \in \mathbb{R}^{k \times n}$$

Principal Components Analysis

- Matrix form

$$\begin{aligned} & \sum_{i=1}^n (\mathbf{w}^T \tilde{\mathbf{x}}_i)^2 \\ &= \sum_{i=1}^n (\mathbf{w}^T \tilde{\mathbf{x}}_i)(\mathbf{w}^T \tilde{\mathbf{x}}_i) \\ &= \sum_{i=1}^n (\mathbf{w}^T \tilde{\mathbf{x}}_i)(\tilde{\mathbf{x}}_i^T \mathbf{w}) \\ &= \mathbf{w}^T \left(\sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \right) \mathbf{w} \\ &= \mathbf{w}^T \tilde{X} \tilde{X}^T \mathbf{w} \end{aligned}$$



Objective function of PCA

$$\max_{W^T W = I} W^T \underbrace{\tilde{X} \tilde{X}^T}_{\text{Covariance matrix}} W$$

Covariance matrix

How to optimize this model?

- Eigen-decomposition for the covariance matrix

$$A = U\Sigma U^T$$

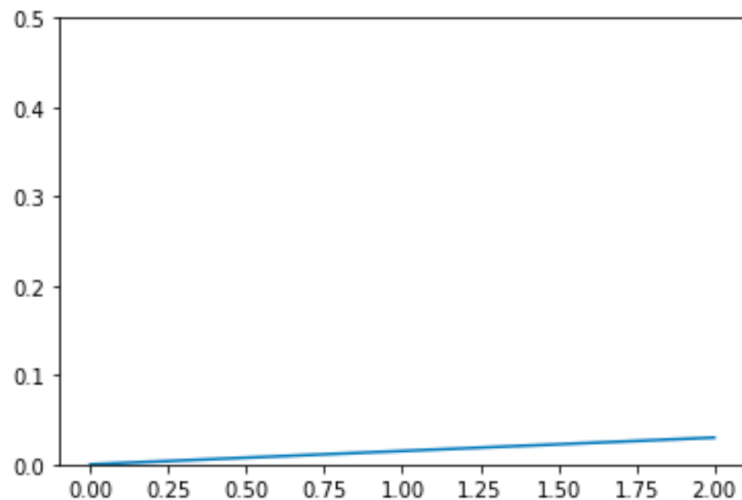
- $A = \tilde{X}\tilde{X}^T \in \mathbb{R}^{d \times d}$ is the covariance matrix
- $U = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d] \in \mathbb{R}^{d \times d}$ where \mathbf{u}_i is the i -th largest eigenvector, $\mathbf{u}_i^T \mathbf{u}_j = 0$, $\|\mathbf{u}_i\|_2 = 1$
- $\Sigma = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_d\}$ where λ_i is the i -th largest eigenvalue $0 \leq \lambda_d \leq \dots \leq \lambda_2 \leq \lambda_1$

- The solution is the largest k eigenvectors

$$\max_{W^T W = I} W^T \tilde{X} \tilde{X}^T W \quad \longrightarrow \quad W = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k] \in \mathbb{R}^{d \times k}$$

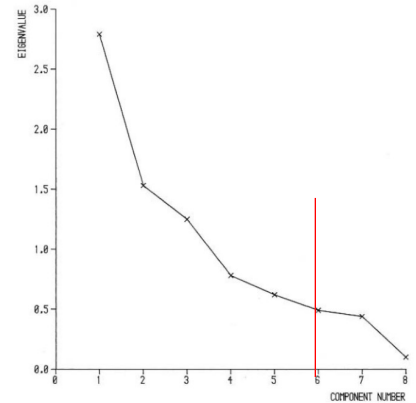
How to optimize this model?

- Interpretation



$$\begin{pmatrix} 2 \\ 0.03 \end{pmatrix} = 2 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 0.03 \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

Covariance matrix $A = U\Sigma U^T$



$$A = \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \lambda_2 \mathbf{u}_2 \mathbf{u}_2^T + \dots + \lambda_k \mathbf{u}_k \mathbf{u}_k^T + \dots + \lambda_d \mathbf{u}_d \mathbf{u}_d^T$$

Only keep the largest k eigenvalue $0 \leq \lambda_d \leq \dots \leq \lambda_2 \leq \lambda_1$

$$A \approx \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \lambda_2 \mathbf{u}_2 \mathbf{u}_2^T + \dots + \lambda_k \mathbf{u}_k \mathbf{u}_k^T$$



$$W = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k] \in \mathbb{R}^{d \times k}$$

Summary

- Step 1: Mean subtraction

$$\tilde{X} = X - \frac{1}{n}X\mathbf{1}\mathbf{1}^T$$

- Step 2: Compute the covariance matrix

$$A = \tilde{X}\tilde{X}^T$$

- Step 3: Eigen-decomposition

$$A = U\Sigma U^T$$

- Step 4: Keep the largest k eigenvectors

$$W = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k] \in \mathbb{R}^{d \times k}$$