Machine Learning

Unsupervised Learning

FAST DISCOVERING THE FUTURE

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● つくで

- ✓ Ingredients of Machine Learning
- ✓ Classification Basics, Basic Linear Classifier
- ✓ K-Nearest Neighbours and Naive Bayes Classifier
- $\checkmark\,$ Linear and Quadratic Discriminant Analysis
- ✓ Support Vector Machines (SVM)
- ✓ Decision Trees
- ✓ Ensemble Methods (Bagging, Weighted Voting, Stacking)
- $\checkmark\,$ Regression Methods
- $\checkmark\,$ Evaluation and Scoring of Classifiers
- ✓ Ensemble Methods (Boosting)
- ✓ Clustering (Hierarchical, K-Means, K-Medoids, DBSCAN)

• • • • • • • • • • •

- Dimensionality Reduction
- Principal Component Analysis (PCA)
- t-Distributed Stochastic Neighbor Embedding (t-SNE)

イロト イ団ト イヨト イヨ

• Working directly with high-dimensional data, such as images, comes with some difficulties

イロト イヨト イヨト イヨ

- Working directly with high-dimensional data, such as images, comes with some difficulties
- High-dimensional data is often overcomplete, i.e., many dimensions are redundant and can be explained by a combination of other dimensions.

Image: A math the second se

- Working directly with high-dimensional data, such as images, comes with some difficulties
- High-dimensional data is often overcomplete, i.e., many dimensions are redundant and can be explained by a combination of other dimensions.
- Dimensions in high-dimensional data are often correlated so that the data possesses an intrinsic lower-dimensional structure.

• • • • • • • • • • • •

- Working directly with high-dimensional data, such as images, comes with some difficulties
- High-dimensional data is often overcomplete, i.e., many dimensions are redundant and can be explained by a combination of other dimensions.
- Dimensions in high-dimensional data are often correlated so that the data possesses an intrinsic lower-dimensional structure.
- Dimensionality reduction exploits structure and correlation and allows us to work with a more compact representation of the data, ideally without losing much information.

(日)

- Working directly with high-dimensional data, such as images, comes with some difficulties
- High-dimensional data is often overcomplete, i.e., many dimensions are redundant and can be explained by a combination of other dimensions.
- Dimensions in high-dimensional data are often correlated so that the data possesses an intrinsic lower-dimensional structure.
- Dimensionality reduction exploits structure and correlation and allows us to work with a more compact representation of the data, ideally without losing much information.
- It can also be useful to detect potential patterns in high dimensial data, by visualizing the first 2-3 projections.

イロト イヨト イヨト イヨト



メロト メタト メヨト メヨ



イロト イヨト イヨト イヨ



メロト メタト メヨト メヨ



What is the problem with **high-dimensional** things?

Hard to visualise



What is the problem with **high-dimensional** things?

Hard to visualise























FAST Foundation





















What is the curse of dimensionality?





On average 55.5% of cells will be either empty or singletons





What is the **curse** of **dimensionality**? (part **II**)





What is the **curse** of **dimensionality**? (part **II**)



Distances become similar in high-dimensional space







Feature extraction vs feature elimination



Feature extraction vs feature elimination



Feature extraction vs feature elimination


Feature extraction vs feature elimination



Feature extraction vs feature elimination



Feature extraction vs feature elimination



Principle Component Analysis



Principle Component #2 (PC2)

Principle Component Analysis



Principle Component #2 (PC2)

1-Dimensional data





X

2-Dimensional data



3-Dimensional data





200-Dimensional data?

FAST Foundation



200-Dimensional data?

Are all of these dimensions equally useful?



























How about we make new axes from these lines?



2-D example revisited

How about we make **new axes** from **these lines**? y y Data is mostly spread along this 6 line 5 4 З And a little bit along 2 this line 1 1 2 З 4 5 X



2-D example revisited

How about we make **new axes** from **these lines**?





These new axes are called principle components



These new axes are called principle components



These new axes are called principle components












Principle components are **not additional axes/dimensions**

How many PCs will be formed in 200D space?



Principle components are **not additional axes/dimensions**







No exceptions, 200 PCs

Principle components are not additional axes/dimensions

How many PCs will be formed in 200D space?





No exceptions, **200 PCs** But what is **the benefit** of having PCs?









Principle components are not additional axes/dimensions

How many PCs will be formed in 200D space?





No exceptions, **200 PCs** But what is **the benefit** of having PCs?

Principle components are not additional axes/dimensions

How many PCs will be formed in 200D space?





First few PCs would be enough to capture important information







x - x	y - ÿ	Î
-2	-2	
-1	0	
0	1	
1	0	
2	1	









What are the **dimensions** of the transposed matrix?









$Z^\top \times Z = S$





2 1





2 1





How to interpret values in covariance matrix?





Covariance matrix

How to interpret values in covariance matrix?





Covariance matrix



















How to **interpret** values in **covariance** matrix




















$$[-2, -1, 0, 1, 2]$$
 $\bar{\mathbf{x}} = 0$ $\sigma = 2.5$





Covariance matrix







Covariance matrix

















How to interpret values in covariance matrix?





Covariance matrix



How to interpret values in covariance matrix?



Covariance indicates how two variables are related. A **positive** covariance means the variables are **positively related**, while a **negative** covariance means the variables are **inversely related**.













How to interpret values in covariance matrix?





Covariance matrix















Covariance **0** means that there is **no relationship** between two variables. Knowing something about the value of one does not say anything about the value of the other.





Covariance matrix








 PDP^{T}







Eigendecomposition



For an example: https://www.scss.tcd.ie/Rozenn.Dahyot/CS1BA1/SolutionEigen.pdf

Dimensionality Reduction

Eigendecomposition





-1 0 1 2 X



-2 -1 0 1 2

-3

X























































Do you still remember what was it all about?



Do you still remember what was it all about? We want to **reduce the dimensionality**!


Do you still remember what was it all about? We want to **reduce the dimensionality**!



We can **ignore** the **second eigenvector** because it does not contain much information



























We can **ignore** the **second PC** because it does not contain much information

How much information the second PC contains?



We can **ignore** the **second PC** because it explains only **10.5%** of variation







No exceptions, 200 PCs



No exceptions, 200 PCs





Variance explained is a good criteria for choosing the total number of PCs to keep

You should keep as many PCs as it takes to explain **90%** of **total variance**



Variance explained is a good criteria for choosing the total number of PCs to keep

You should keep as many PCs as it takes to explain **90%** of **total variance**



Principle Component Analysis (**PCA**)



Can be used as part of supervised learning pipeline

E			
	 αun	dat	ion
1/10	oun	uau	

Dimensionality Reduction























PCA has an "undo" button



You can recover the original features back!

EAST Foundati	on
	<u>u</u>
PCA has an "undo" button

Conventional PCA $eigenvectors \times Z_{original}^{\mathsf{T}} = Z_{transformed}$

Reversed PCA
eigenvectors^T ×
$$Z_{transformed} = Z_{original}$$

You can recover the original features back!

Convince yourself: http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf

FAST Foundation

Dimensionality Reduction

Principal Component Analysis

Principle components are linear combinations of original features



FAST Foundation

Principle components are linear combinations of original features



So if you predict anything based on PCs, the **meaning** of original features **is not preserved** after the transformation

• In PCA, we are interested in finding projections \mathbf{x}'_n of data points \mathbf{x}_n that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.

イロト イヨト イヨト イヨ

- In PCA, we are interested in finding projections \mathbf{x}'_n of data points \mathbf{x}_n that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.
- Suppose we have a dataset $\mathbf{X} = {\mathbf{x}_1, \dots, \mathbf{x}_N} \in \mathbb{R}^{N \times D}$, where $\mathbf{x}_n \in \mathbb{R}^D$, $\mathbb{E}(\mathbf{X}) = 0$ and the sample covariance of \mathbf{X} is

$$S = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^T$$

イロト イヨト イヨト イヨ

- In PCA, we are interested in finding projections \mathbf{x}'_n of data points \mathbf{x}_n that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.
- Suppose we have a dataset $\mathbf{X} = {\mathbf{x}_1, \dots, \mathbf{x}_N} \in \mathbb{R}^{N \times D}$, where $\mathbf{x}_n \in \mathbb{R}^D$, $\mathbb{E}(\mathbf{X}) = 0$ and the sample covariance of \mathbf{X} is

$$S = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^T$$

• We assume there exists a low-dimensional compressed representation (code) of \mathbf{x}_n

$$\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n \in \mathbb{R}^M,$$

where $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$ is the projection matrix.

・ロト ・日 ト ・ヨト ・ヨト

- In PCA, we are interested in finding projections x'_n of data points x_n that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.
- Suppose we have a dataset $\mathbf{X} = {\mathbf{x}_1, \dots, \mathbf{x}_N} \in \mathbb{R}^{N \times D}$, where $\mathbf{x}_n \in \mathbb{R}^D$, $\mathbb{E}(\mathbf{X}) = 0$ and the sample covariance of \mathbf{X} is

$$S = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^T$$

• We assume there exists a low-dimensional compressed representation (code) of \mathbf{x}_n

$$\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n \in \mathbb{R}^M,$$

where $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$ is the projection matrix.

• We further assume that the columns of **B** are orthonormal, so that $\mathbf{b}_i^T \mathbf{b}_j = 0$ iff $i \neq j$ and $\mathbf{b}_i^T \mathbf{b}_i = 1$

・ロト ・ 日 ト ・ 日 ト ・ 日 ト ・

- In PCA, we are interested in finding projections \mathbf{x}'_n of data points \mathbf{x}_n that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.
- Suppose we have a dataset $\mathbf{X} = {\mathbf{x}_1, \dots, \mathbf{x}_N} \in \mathbb{R}^{N \times D}$, where $\mathbf{x}_n \in \mathbb{R}^D$, $\mathbb{E}(\mathbf{X}) = 0$ and the sample covariance of \mathbf{X} is

$$S = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^T$$

• We assume there exists a low-dimensional compressed representation (code) of \mathbf{x}_n

$$\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n \in \mathbb{R}^M,$$

where $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$ is the projection matrix.

- We further assume that the columns of B are orthonormal, so that $\mathbf{b}_i^T \mathbf{b}_j = 0$ iff $i \neq j$ and $\mathbf{b}_i^T \mathbf{b}_i = 1$
- We seek an $M\text{-dimensional subspace }U\subseteq \mathbb{R}^D, \, dim(U)=M < D$ onto which we project the data

・ロト ・四ト ・ヨト ・ヨト

- In PCA, we are interested in finding projections \mathbf{x}'_n of data points \mathbf{x}_n that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.
- Suppose we have a dataset $\mathbf{X} = {\mathbf{x}_1, \dots, \mathbf{x}_N} \in \mathbb{R}^{N \times D}$, where $\mathbf{x}_n \in \mathbb{R}^D$, $\mathbb{E}(\mathbf{X}) = 0$ and the sample covariance of \mathbf{X} is

$$S = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^T$$

• We assume there exists a low-dimensional compressed representation (code) of \mathbf{x}_n

$$\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n \in \mathbb{R}^M,$$

where $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$ is the projection matrix.

- We further assume that the columns of B are orthonormal, so that $\mathbf{b}_i^T \mathbf{b}_j = 0$ iff $i \neq j$ and $\mathbf{b}_i^T \mathbf{b}_i = 1$
- We seek an $M\text{-dimensional subspace }U\subseteq \mathbb{R}^D, \, dim(U)=M < D$ onto which we project the data
- We denote the projected data as $\mathbf{x}'_n \in U$ and their coordinates by \mathbf{z}_n

• We can describe the information contained in the data by looking at the spread of the data and measure it with **variance**

イロト イ団ト イヨト イヨ

- We can describe the information contained in the data by looking at the spread of the data and measure it with **variance**
- In this setting, retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional representation

・ロト ・ 日 ・ ・ ヨ ト ・

- We can describe the information contained in the data by looking at the spread of the data and measure it with **variance**
- In this setting, retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional representation
- We maximize the variance of the low-dimensional code using a sequential approach.

• • • • • • • • • • • • •

- We can describe the information contained in the data by looking at the spread of the data and measure it with variance
- In this setting, retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional representation
- We maximize the variance of the low-dimensional code using a sequential approach.
- We start by seeking a single vector $\mathbf{b}_1 \in \mathbb{R}^D$ that maximizes variance of the projected data

$$\mathbb{V}ar(z_1) = \frac{1}{N} \sum_{n=1}^{N} z_{1n}^2 = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{b}_1^T \mathbf{x}_n)^2 =$$

- We can describe the information contained in the data by looking at the spread of the data and measure it with **variance**
- In this setting, retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional representation
- We maximize the variance of the low-dimensional code using a sequential approach.
- We start by seeking a single vector $\mathbf{b}_1 \in \mathbb{R}^D$ that maximizes variance of the projected data

$$\mathbb{V}ar(z_1) = \frac{1}{N} \sum_{n=1}^{N} z_{1n}^2 = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{b}_1^T \mathbf{x}_n)^2 = \\ = \frac{1}{N} \sum_{n=1}^{N} \mathbf{b}_1^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{b}_1 = \mathbf{b}_1^T (\frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^T) \mathbf{b}_1 = \mathbf{b}_1^T S \mathbf{b}_1,$$

where S is the sample covariance matrix.

(日)

- We can describe the information contained in the data by looking at the spread of the data and measure it with **variance**
- In this setting, retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional representation
- We maximize the variance of the low-dimensional code using a sequential approach.
- We start by seeking a single vector $\mathbf{b}_1 \in \mathbb{R}^D$ that maximizes variance of the projected data

$$\mathbb{V}ar(z_1) = \frac{1}{N} \sum_{n=1}^{N} z_{1n}^2 = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{b}_1^T \mathbf{x}_n)^2 = \\ = \frac{1}{N} \sum_{n=1}^{N} \mathbf{b}_1^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{b}_1 = \mathbf{b}_1^T (\frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^T) \mathbf{b}_1 = \mathbf{b}_1^T S \mathbf{b}_1,$$

where S is the sample covariance matrix.

Note that if we would not put a unit vector contraint on b₁, the variance would increase for longer-length b₁.

We need to solve the following constrained optimization problem

 $\max_{\mathbf{b}_1} \mathbf{b}_1^T S \mathbf{b}_1$

subject to $\|\mathbf{b}_1\|^2 = 1$

イロト イヨト イヨト イヨ

We need to solve the following constrained optimization problem

 $\max_{\mathbf{b}_1} \mathbf{b}_1^T S \mathbf{b}_1$

subject to $\|\mathbf{b}_1\|^2 = 1$

The Lagragian will be

$$\Lambda(\mathbf{b}_1, \lambda_1) = \mathbf{b}_1^T S \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^T \mathbf{b}_1)$$

イロト イヨト イヨト イヨ

We need to solve the following constrained optimization problem

$$\max_{\mathbf{b}_1} \mathbf{b}_1^T S \mathbf{b}_1$$

subject to $\|\mathbf{b}_1\|^2 = 1$

The Lagragian will be

$$\Lambda(\mathbf{b}_1, \lambda_1) = \mathbf{b}_1^T S \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^T \mathbf{b}_1)$$

Finding the partial derivatives w.r.t. \mathbf{b}_1 and λ_1

$$\frac{\partial \Lambda}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^T S - 2\lambda_1 \mathbf{b}_1^T, \ \frac{\partial \Lambda}{\partial \lambda_1} = 1 - \mathbf{b}_1^T \mathbf{b}_1$$

イロト イロト イヨト イヨ

We need to solve the following constrained optimization problem

 $\max_{\mathbf{b}_1} \mathbf{b}_1^T S \mathbf{b}_1$

subject to $\|\mathbf{b}_1\|^2 = 1$

The Lagragian will be

$$\Lambda(\mathbf{b}_1, \lambda_1) = \mathbf{b}_1^T S \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^T \mathbf{b}_1)$$

Finding the partial derivatives w.r.t. \mathbf{b}_1 and λ_1

$$\frac{\partial \Lambda}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^T S - 2\lambda_1 \mathbf{b}_1^T, \ \frac{\partial \Lambda}{\partial \lambda_1} = 1 - \mathbf{b}_1^T \mathbf{b}_1$$

Setting these partial derivatives to $0 \ {\rm gives} \ {\rm us}$

$$S\mathbf{b}_1 = \lambda_1 \mathbf{b}_1$$
$$\mathbf{b}_1^T \mathbf{b}_1 = 1$$

Does this look familiar?

FAST Foundation

イロト イヨト イヨト イヨ

 $S\mathbf{b}_1 = \lambda_1 \mathbf{b}_1$ $\mathbf{b}_1^T \mathbf{b}_1 = 1$

 \mathbf{b}_1 is an eigenvector of the data covariance matrix S, and the Lagrange multiplier λ_1 plays the role of the corresponding eigenvalue.

イロト イヨト イヨト イヨ

$$S\mathbf{b}_1 = \lambda_1 \mathbf{b}_1$$
$$\mathbf{b}_1^T \mathbf{b}_1 = 1$$

 \mathbf{b}_1 is an eigenvector of the data covariance matrix S, and the Lagrange multiplier λ_1 plays the role of the corresponding eigenvalue. The objective can be re-written as

$$\mathbb{V}ar(z_1) = \mathbf{b}_1^T S \mathbf{b}_1 = \lambda_1 \mathbf{b}_1^T \mathbf{b}_1 = \lambda_1,$$

the variance of the data projected onto a one-dimensional subspace equals the eigenvalue that is associated with the basis vector \mathbf{b}_1 that spans this subspace.

< □ > < 同 > < 回 > < Ξ > < Ξ

$$S\mathbf{b}_1 = \lambda_1 \mathbf{b}_1$$
$$\mathbf{b}_1^T \mathbf{b}_1 = 1$$

 \mathbf{b}_1 is an eigenvector of the data covariance matrix S, and the Lagrange multiplier λ_1 plays the role of the corresponding eigenvalue. The objective can be re-written as

$$\mathbb{V}ar(z_1) = \mathbf{b}_1^T S \mathbf{b}_1 = \lambda_1 \mathbf{b}_1^T \mathbf{b}_1 = \lambda_1,$$

the variance of the data projected onto a one-dimensional subspace equals the eigenvalue that is associated with the basis vector \mathbf{b}_1 that spans this subspace. To obtain the projection of \mathbf{x}_n on the obtained subspace, we can use the following formula

$$\mathbf{x}_n' = \mathbf{b}_1 z_{1n} = \mathbf{b}_1 \mathbf{b}_1^T \mathbf{x}_n \in \mathbb{R}^D$$

イロト イ団ト イヨト イヨト

• Assume we have found the first m-1 principal components as the m-1 eigenvectors of S that are associated with the largest m-1 eigenvalues.

・ロト ・日下・ ・ ヨト・

- Assume we have found the first m-1 principal components as the m-1 eigenvectors of S that are associated with the largest m-1 eigenvalues.
- The *m*-th principal component can be found by subtracting the effect of the first m-1 principal components $\mathbf{b}_1, \ldots, \mathbf{b}_{m-1}$ from the data, by trying to find principal components that compress the remaining information.

$$\hat{\mathbf{X}} = \mathbf{X} - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T \mathbf{X} = \mathbf{X} - \mathbf{B}_{m-1} \mathbf{X},$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ and $\mathbf{B}_{m-1} = \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T$ is the projection matrix onto the subspace spanned by $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$.

< □ > < □ > < □ > < □ > < □ >

- Assume we have found the first m-1 principal components as the m-1 eigenvectors of S that are associated with the largest m-1 eigenvalues.
- The *m*-th principal component can be found by subtracting the effect of the first m-1 principal components $\mathbf{b}_1, \ldots, \mathbf{b}_{m-1}$ from the data, by trying to find principal components that compress the remaining information.

$$\hat{\mathbf{X}} = \mathbf{X} - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T \mathbf{X} = \mathbf{X} - \mathbf{B}_{m-1} \mathbf{X},$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ and $\mathbf{B}_{m-1} = \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T$ is the projection matrix onto the subspace spanned by $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$.

• The m-th PC can be found by maximizing the variance

$$\begin{split} \mathbb{V}ar(z_m) &= \frac{1}{N}\sum_{n=1}^N z_{mn}^2 = \frac{1}{N}\sum_{n=1}^N (\mathbf{b}_m^T\hat{\mathbf{x}_n})^2 = \mathbf{b}_m^T\hat{S}\mathbf{b}_m\\ \text{subject to } \|\mathbf{b}_m\|^2 = 1\\ \hat{S} \text{ is the covariance matrix of } \hat{\mathbf{X}} \end{split}$$

where

< □ > < □ > < □ > < □ > < □ >

- Assume we have found the first m-1 principal components as the m-1 eigenvectors of S that are associated with the largest m-1 eigenvalues.
- The *m*-th principal component can be found by subtracting the effect of the first m-1 principal components $\mathbf{b}_1, \ldots, \mathbf{b}_{m-1}$ from the data, by trying to find principal components that compress the remaining information.

$$\hat{\mathbf{X}} = \mathbf{X} - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T \mathbf{X} = \mathbf{X} - \mathbf{B}_{m-1} \mathbf{X},$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ and $\mathbf{B}_{m-1} = \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T$ is the projection matrix onto the subspace spanned by $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$.

• The *m*-th PC can be found by maximizing the variance

$$\mathbb{V}ar(z_m) = \frac{1}{N} \sum_{n=1}^{N} z_{mn}^2 = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{b}_m^T \hat{\mathbf{x}}_n)^2 = \mathbf{b}_m^T \hat{S} \mathbf{b}_m$$

subject to $\|\mathbf{b}_m\|^2 = 1$
where \hat{S} is the covariance matrix of $\hat{\mathbf{X}}$
• \mathbf{b}_m is also an eigenvector of S and $\mathbb{V}ar(z_m) = \mathbf{b}_m^T S \mathbf{b}_m = \lambda_m$ is the *m*-th
largest eigenvalue of S

FAST Foundation

• We can reduce the dimensionality of our data, by using the M eigenvectors of the covariance matrix S associated with the M largest eigenvalues.

イロト イヨト イヨト イヨ

- We can reduce the dimensionality of our data, by using the M eigenvectors of the covariance matrix S associated with the M largest eigenvalues.
- ${\ensuremath{\bullet}}$ The projection matrix ${\ensuremath{\mathbf{B}}}$ consists of the eigenvectors of S
- The amount of variance PCA captured with the first M principal components is



< □ > < 同 > < 回 > < Ξ > < Ξ

- We can reduce the dimensionality of our data, by using the M eigenvectors of the covariance matrix S associated with the M largest eigenvalues.
- ${\ensuremath{\bullet}}$ The projection matrix ${\ensuremath{\mathbf{B}}}$ consists of the eigenvectors of S
- The amount of variance PCA captured with the first M principal components is



• The variance lost by data compression via PCA is

$$\sum_{i=M+1}^D \lambda_m$$

< □ > < 同 > < 回 > < Ξ > < Ξ



・ロト ・ 日 ト ・ 日 ト ・ 日



イロト イヨト イヨト イヨト



イロト イヨト イヨト イヨト



メロト メタト メヨト メヨト



イロト イヨト イヨト イヨト



イロト イヨト イヨト イヨト

t-SNE iteratively tries to make distances in lowdimensional space to be similar to distances in highdimensional space



イロト イ団ト イヨト イヨ
- ✓ Dimensionality Reduction
- ✓ Principal Component Analysis
- ✓ t-Distributed Stochastic Neighbor Embedding (t-SNE)

Image: A math the second se