

Assignment Project Exam Help

Add WeChat edu\_assist\_pro

Dimensionality Reduction with  
Principal <sup>Assignment Project Exam Help</sup> Analysis  
<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Liang Zheng

Australian National University

liang.zheng@anu.edu.au

Meta Sim: Learning to Generate Synthetic Datasets. Kar et al., ICCV 2019

Add WeChat [edu\\_assist\\_pro](#)

[Assignment](#) [Project](#) [Exam](#) [Help](#)

<https://eduassistpro.github.io/>

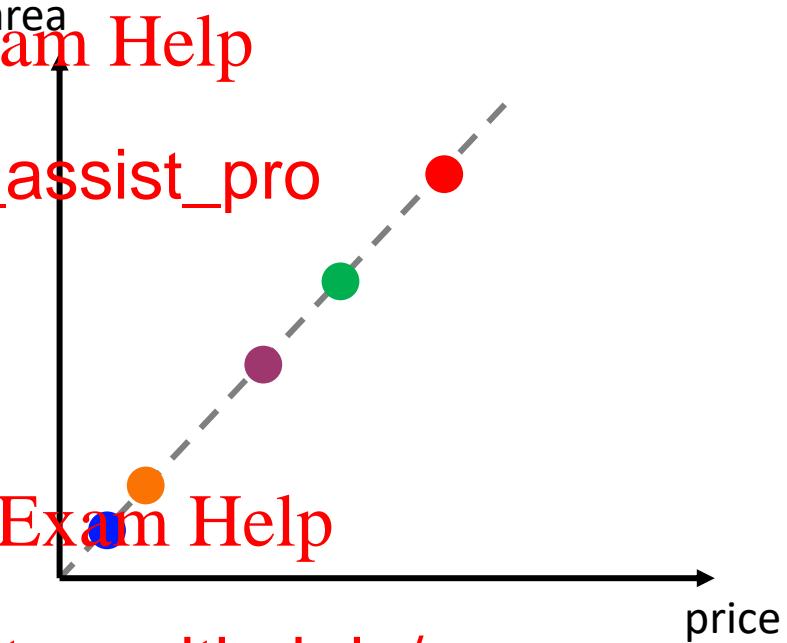
Add WeChat [edu\\_assist\\_pro](#)

# Idea of PCA

Assignment Project Exam Help

	House price (million)	House area (100)
a	10	10
b	2	2
c	7	7
d	1	1
e	5	

Assignment Project Exam Help

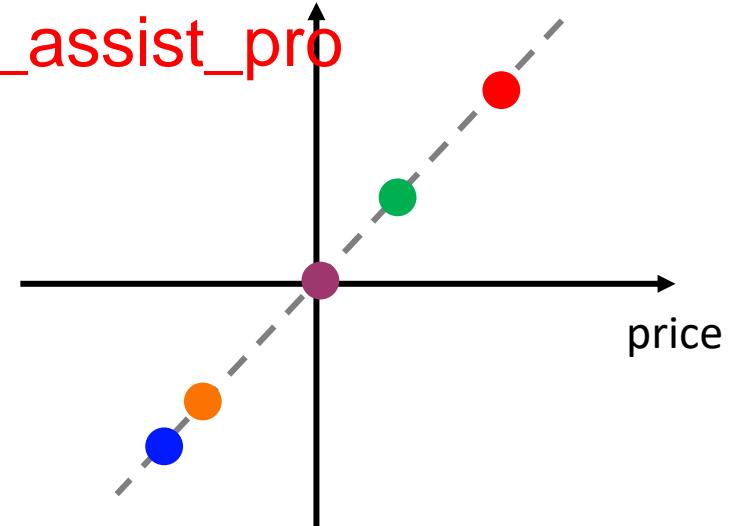


<https://eduassistpro.github.io/>

We subtract means from data points

Assignment Project Exam Help

	House price (normalised)	House area (normalised)
a	5	5
b	-3	-3
c	2	2
d	-4	-4
e	0	0



# Idea of PCA

Assignment Project Exam Help

first principal component

area

Add WeChat edu\_assist\_pro

We rotate our  
price and area axis

Assignment Project Exam Help

<https://eduassistpro.github.io/>

House price  
(normalised)

House area  
(normalised)

	House price (normalised)	House area (normalised)
a	5	5
b	-3	-3
c	2	2
d	-4	-4
e	0	0

Add WeChat edu\_assist\_pro

First principal component

Second principal component

a 7.07

0

b -4.24

0

c 2.82

0

d -5.66

0

e 0

0

# Motivation Assignment Project Exam Help

- High-dimensional data, such as images, are expensive to analyze, interpretate, and visualize, and expensive to store.
- Good news
- high-dimensional data is often overcomplete, i.e., many dimensions are redundant and can be explained by a combination of other dimensions
- Furthermore, dimensions in high-dimensional data are often correlated so that the data possess a natural structure.

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

The data in (a) does not vary much in the  $x_2$ -direction, so that we can express it as if it were on a line – with nearly no loss; see (b).

To describe the data in (b), only the  $x_2$ -coordinate is required, and the data lies in a one-dimensional subspace of  $\mathbb{R}^2$ .

# 10.1 Problem Setting Assignment Project Exam Help

- In PCA, we are interested in finding a set of data points  $\mathbf{x}_n$  that are as similar to the original data points  $\mathbf{x}_n$  as possible, but which have a significantly lower intrinsic dimensionality
- We consider an i.i.d. dataset  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, \mathbf{x}_n \in \mathbb{R}^D$ , with mean  $\mathbf{0}$  that possesses the data covariance matrix

$$\Sigma = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top$$

- We assume there exist <https://eduassistpro.github.io/> for a compressed representation (code)

of  $\mathbf{x}_n$ , where we define the projection matrix  $\mathbf{B}$  such that  $\mathbf{z}_n = \mathbf{B}^\top \mathbf{x}_n$

$$\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$$

# Assignment Project Exam Help

- Example (Coordinate Representation)  
Add WeChat edu\_assist\_pro
- Consider  $\mathbb{R}^2$  with the canonical basis  $e_1 = [1, 0]^T$ ,  $e_2 = [0, 1]^T$ .
- $x \in \mathbb{R}^2$  can be represented as a linear combination of these basis vectors, e.g.,

$$\begin{bmatrix} 5 \\ 3 \end{bmatrix} = 5e_1 + 3e_2$$

- However, when we consider vectors of the form

$$\begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

<https://eduassistpro.github.io/>

they can always be writ

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

- To represent these vectors it is sufficient to define the coordinate code  $z$  of  $\tilde{x}$  with respect to the  $e_2$  vector.

# 10.2 PCA from Maximum Variance Perspective

Add WeChat [edu\\_assist\\_pro](#)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

- We ignore  $x_2$ -coordinate of the data because it did not add too much information: the compressed data (b) is similar to the original data in (a)
- We derive PCA so as to maximize the variance in the low-dimensional representation of the data to retain as much information as possible
- Retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional code (Hotelling, 1933)

## 10.2.1 Direction with Maximal Variance

Add WeChat `edu_assist_pro`

- Data centering
- In the data covariance matrix, we assume centered data.

$$S = \frac{1}{N} \sum_{n=1}^N x_n x_n^T$$

Assignment Project Exam Help

- Let us assume that  $\mu$  sing the properties of the variance, we obtain <https://eduassistpro.github.io/>

$$\mathbb{V}_z[z] = \mathbb{V}_x[(\quad)]_x[\quad] = \mathbb{V}_x[B^T x]$$

- That is, the variance of the low-dimensional ~~is not depend on the mean of the data~~ ~~mean of the data~~.

- With this assumption the mean of the low-dimensional code is also  $\mathbf{0}$  since

$$\mathbb{E}_z[z] = \mathbb{E}_x[B^T x] = B^T \mathbb{E}_x[x] = \mathbf{0}$$

# Assignment Project Exam Help

- To maximize the variance of the low-dimensional representation  $\mathbf{z}_1 \in \mathbb{R}^D$  that maximizes the variance of the first

ode, we first seek a single projected data, i.e., we aim to maximize the variance of  $\mathbf{z} \in \mathbb{R}^M$  so that

$$V_1 := \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1n}^2$$

is maximized, where we defined  $z_{1n}$  as the first coordinate of the low-dimensional representation  $\mathbf{z}_n \in \mathbb{R}^M$  of  $\mathbf{x}_n \in \mathbb{R}^D$ .  $z_{1n}$  is given by,

<https://eduassistpro.github.io/>

i.e., it is the coordinate of  $\mathbf{x}_n$  onto the one-dimensional subspace spanned by  $\mathbf{b}_1$ . We

$$\begin{aligned} V_1 &= \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_1^\top \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n^\top \mathbf{b}_1 \\ &= \mathbf{b}_1^\top \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{b}_1 = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \end{aligned}$$

where  $\mathbf{S}$  is the data covariance matrix.

- We further restrict all solutions to  $\|\mathbf{b}_1\|^2 = 1$

# Assignment Project Exam Help

- We have the following constrained optimization problem

$$\begin{array}{ll} \text{Add WeChat} & \max_{\mathbf{b}_1} \\ \text{subject t} & \parallel \mathbf{b}_1 \parallel \end{array}$$

- We obtain the Lagrangian (not required in this course),

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

- The partial derivatives of  $\mathfrak{L}$  with respect to  $\mathbf{b}_1$  and  $\lambda_1$  are

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

- Setting these partial deriv

$$\mathbf{S} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1$$

$$\text{Add WeChat} \text{edu\_assist\_pro}$$

- We see that  $\mathbf{b}_1$  is an eigenvector of  $\mathbf{S}$ , and  $\lambda_1$  is the corresponding eigenvalue. We rewrite our objective as,

$$V_1 = \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1^T \mathbf{b}_1 = \lambda_1$$

- i.e., 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 maximize the variance of the low-dimensional code, we choose the basis vector associated with the largest eigenvalue of the data covariance matrix. This eigenvector is called the first principal component.

## Assignment Project Exam Help

### 10.2.2 $M$ -dimensional Subspace with Maximal Variance

- Assume we have found the  $m$  principal components  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$  that are associated with the largest  $m - 1$  eigenvalues.
- We want to find the  $m$ th principal component.
- We subtract the effect of the first  $m - 1$  principal components  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$  from the data, and find principal components that compress the remaining information. We then arrive at the new data matrix,

$\hat{\mathbf{X}}$  <https://eduassistpro.github.io/>

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

- To find the  $m$ th principal component, we maximize the variance

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

subject to  $\|\mathbf{b}_m\|^2 = 1$ , and we define  $\hat{\mathbf{S}}$  as the data covariance matrix of the transformed dataset  $\hat{\mathbf{X}} := \{\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_N\}$ .

# Assignment Project Exam Help

Add WeChat edu\_assist\_pro

- The optimal solution  $\mathbf{b}_m$  is the eigenvector of  $\mathbf{S}$  that is associated with the largest eigenvalue of  $\widehat{\mathbf{S}}$ .
- In fact, we can derive that

$$\widehat{\mathbf{S}}\mathbf{b}_m = \mathbf{S}\mathbf{b}_m = \lambda_m \mathbf{b}_m \quad (1)$$

- $\mathbf{b}_m$  is not only an eigenvector of  $\mathbf{S}$  but also of  $\widehat{\mathbf{S}}$ .
- Specifically,  $\lambda_m$  is the largest eigenvalue of  $\mathbf{S}$  and  $\lambda_m$  is the  $m$ th largest eigenvalue of  $\widehat{\mathbf{S}}$ , and both have the associated eigenvector  $\mathbf{b}_m$ .
- Moreover,  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$  are associated with eigenvalue 0.
- Considering (1) and,  $\mathbf{b}_m^T \mathbf{b}_m = 1$ , the variance of the data projected onto the  $m$ th principal component is

$$V_m = \mathbf{b}_m^T \mathbf{S} \mathbf{b}_m = \lambda_m \mathbf{b}_m^T \mathbf{b}_m = \lambda_m$$

- This means that the variance of the data, when projected onto an  $M$ -dimensional subspace, equals the sum of the eigenvalues that are associated with the corresponding eigenvectors of the data covariance matrix.

# MNIST dataset

Assignment Project Exam Help

- 60,000 examples of handwritten digit
- Each digit is a grayscale image of size  $28 \times 28$ , i.e., it contains 784 pixels.
- We can interpret every image in this dataset as a vector  $x \in \mathbb{R}^{784}$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Example - Eigenvalues of MNIST digit “8”

Assignment Project Exam Help  
Add WeChat [edu\\_assist\\_pro](#)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

(a) Top 200 largest eigenvalues captured by the components.

Add WeChat [edu\\_assist\\_pro](#)

- A 784-dim vector is used to represent an image
- Taking all images of “8” in MNIST, we compute the eigenvalues of the data covariance matrix.
- We see that only a few of them have a value that differs significantly from 0.
- Most of the variance, when projecting data onto the subspace spanned by the corresponding eigenvectors, is captured by only a few principal components

Overall

# Assignment Project Exam Help

Add WeChat `edu_assist_pro`

- To find an  $M$ -dimensional subspace of  $\mathbb{R}^D$  that retains as much information as possible,
- We choose the columns of  $B = [b_1, \dots, b_M] \in \mathbb{R}^{D \times M}$  as the  $M$  eigenvectors of the data covariance matrix  $S$  that are associated with the  $M$  largest eigenvalues.
- The maximum amount of variance retained by the first  $M$  principal components is

<https://eduassistpro.github.io/>

Add WeChat  $\sum_{m=1}^{V_M} \lambda_m$  `edu_assist_pro`

where the  $\lambda_m$  are the  $M$  largest eigenvalues of the data covariance matrix  $S$ .

- The variance lost by data compression via PCA is

$$J_M = \sum_{j=M+1}^D \lambda_j = V_D - V_M$$

- Instead of these absolute quantities, we can define the relative variance captured as  $\frac{V_M}{V_D}$ , and the relative variance lost by compression as  $1 - \frac{V_M}{V_D}$ .

## 10.3 PCA from Projection Perspective

Add WeChat `edu_assist_pro`

- Previously, we derived PCA by maximizing the variance in the projected space to retain as much information as possible

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

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

Assignment Project Exam Help

- Alternatively, we derive PCA as an algorithm that directly minimizes the average reconstruction

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

## 10.3.1 Setting and Objective

Add WeChat [edu\\_assist\\_pro](#)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

(a) A vector  $\mathbf{x} \in \mathbb{R}^2$  (red cross) shall be projected onto a one-dimensional subspace  $U \subseteq \mathbb{R}^2$  spanned by  $\mathbf{b}$ . The resulting vector  $\tilde{\mathbf{x}}$  is shown by the red line segment connecting  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$ .

- We wish to project  $\mathbf{x}$  to  $\tilde{\mathbf{x}}$  in a lower-dimensional space, such that  $\tilde{\mathbf{x}}$  is similar to the original data point  $\mathbf{x}$ . That is,
- We aim to minimize the (Euclidean) distance  $\|\mathbf{x} - \tilde{\mathbf{x}}\|$

# Assignment Project Exam Help

Add WeChat edu\_assist\_pro

- Given an orthonormal basis  $(\mathbf{b}_1, \dots, \mathbf{b}_D)$  of  $\mathbb{R}^D$ , any  $\mathbf{x} \in \mathbb{R}^D$  can be written as a linear combination of the basis vectors of  $\mathbb{R}^D$ :

$$\mathbf{x} = \sum_{d=1}^D \zeta_d \mathbf{b}_d = \sum_{m=1}^M \zeta_m \mathbf{b}_m + \sum_{j=M+1}^D \zeta_j \mathbf{b}_j$$

Assignment Project Exam Help

for suitable coordinates

- We aim to find vectors  $\mathbf{U} \subseteq \mathbb{R}^D$ ,  $\dim(\mathbf{U}) = M$ , so that

Add WeChat edu\_assist\_pro

$$\tilde{\mathbf{x}} = \sum_{m=1}^M z_m \mathbf{b}_m$$

is as similar to  $\mathbf{x}$  as possible.

# Assignment Project Exam Help

- We have a dataset  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, \mathbf{x}_n \in \mathbb{R}^D$  centered at  $\mathbf{0}$ , i.e.,  $\mathbb{E}[\mathbf{x}] = \mathbf{0}$ .
- We want to find the best linear projection of  $\mathcal{X}$  onto a lower dimensional subspace  $\mathbf{U} \subseteq \mathbb{R}^D$ ,  $\dim(\mathbf{U}) = M$ . Also, we want to find an orthonormal basis vectors  $\mathbf{b}_1, \dots, \mathbf{b}_M$ .
- We call this subspace  $\mathbf{U}$  the principal subspace.
- The projections of the data points are denoted by

Assignment Project Exam Help

<https://eduassistpro.github.io/>

where  $\mathbf{z}_n := [z_{1n}, \dots, z_{Mn}] \in \mathbb{R}^M$  is the coordinate vector of  $\tilde{\mathbf{x}}_n$  with respect to the basis  $(\mathbf{b}_1, \dots, \mathbf{b}_M)$ .

- We want to have  $\tilde{\mathbf{x}}_n$  as similar to  $\mathbf{x}_n$  as possible.
- We define our objective as minimizing the average squared Euclidean distance (reconstruction error)

$$J_M := \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$$

- We need to find the orthonormal basis of the principal subspace and the coordinates  $\mathbf{z}_n \in \mathbb{R}^M$  of the projections with respect to this basis.

Assignment Project Exam Help

## 10.3.2 Finding Optimal Coordinates

Add WeChat [edu\\_assist\\_pro](#)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

(a) A vector  $\mathbf{x} \in \mathbb{R}^2$  (red cross) shall be projected onto a one-dimensional subspace  $U \subseteq \mathbb{R}^2$  spanned by  $\mathbf{b}$ . The red lines represent the vectors  $\mathbf{x} - \tilde{\mathbf{x}}_i$  for 50 different  $\tilde{\mathbf{x}}_i$ .

- We want to find  $\tilde{\mathbf{x}}$  in a subspace spanned by  $\mathbf{b}$  that minimizes  $\|\mathbf{x} - \tilde{\mathbf{x}}\|$ .
- Apparently, this will be the orthogonal projection

# Assignment Project Exam Help

Add WeChat edu\_assist\_pro

$$J_M := \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2$$

- Given an ONB  $(b_1, \dots, b_M)$  of  $U \subseteq \mathbb{R}^D$ , to find the optimal coordinates  $z_m$  with respect to this basis, we calculate the partial derivatives

$$\frac{\partial J_M}{\partial z_{in}} = \frac{\partial J_M}{\partial \tilde{x}_n} \frac{\partial \tilde{x}_n}{\partial z_{in}}$$

Assignment Project Exam Help

$$\frac{\partial}{\partial z_{in}} \left( \frac{\partial J_M}{\partial \tilde{x}_n} \right) = \frac{\partial}{\partial z_{in}} \left( \frac{2}{N} (x_n - \tilde{x}_n)^T \right)$$

$$\tilde{x}_n := \sum_{m=1}^M z_{mn} b_m = B z_n \in \mathbb{R}^D$$

for  $i = 1, \dots, M$ , such that we obtain

$$\frac{\partial J_M}{\partial z_{in}} = -\frac{2}{N} (x_n - \tilde{x}_n)^T b_i = -\frac{2}{N} \left( x_n - \sum_{m=1}^M z_{mn} b_m \right)^T b_i$$

ONB

$$= -\frac{2}{N} (x_n^T b_i - z_{in} b_i^T b_i) = -\frac{2}{N} (x_n^T b_i - z_{in})$$

# Assignment Project Exam Help

Add  $\frac{\partial J_M}{\partial z_{in}}$   $\frac{2}{N}$  WeChat [edu\\_assist\\_pro](#)

- Setting this partial derivative to 0 yields immediately the optimal coordinates

$$z_{in} = \mathbf{x}_n^T \mathbf{b}_i = \mathbf{b}_i^T \mathbf{x}_n$$

for  $i = 1, \dots, M$ , and  $n = 1, \dots, N$ .

# Assignment Project Exam Help

- The optimal coordinates  $z_{in}$  of the projection  $\tilde{\mathbf{x}}_n$  are the coordinates of the orthogonal projection onto the one-dimensional subspace that is spanned by <https://eduassistpro.github.io/>
- The optimal linear projection  $\tilde{\mathbf{x}}_n$  of  $\mathbf{x}_n$  is the principal component projection.
- The coordinates of  $\tilde{\mathbf{x}}_n$  with respect to the basis  $(\mathbf{b}_1, \dots, \mathbf{b}_M)$  are the coordinates of the orthogonal projection of  $\mathbf{x}_n$  onto the principal subspace.

Add WeChat [edu\\_assist\\_pro](#)

# Assignment Project Exam Help

Add WeChat `edu_assist_pro`

## Assignment Project Exam Help

- (a) A vector  $\mathbf{x} \in \mathbb{R}^2$  (re projected onto a one-d  $U \subseteq \mathbb{R}^2$  spanned by  $\mathbf{b}$ )
- Distances  $\|\mathbf{x} - \tilde{\mathbf{x}}_i\|$  for 50 different  $\tilde{\mathbf{x}}_i$  by the red lines
- <https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

- (c) Distances  $\|\mathbf{x} - \tilde{\mathbf{x}}\|$  for some  $\tilde{\mathbf{x}} = z_1 \mathbf{b} \in U = \text{span}[\mathbf{b}]$

- (d) The vector  $\tilde{\mathbf{x}}$  that minimizes  $\|\mathbf{x} - \tilde{\mathbf{x}}\|$  is the orthogonal projection of  $\mathbf{x}$  onto  $U$ .

# Assignment Project Exam Help

- We briefly recap orthogonal projections 3.8 (Analytic geometry).
- If  $(\mathbf{b}_1, \dots, \mathbf{b}_D)$  is an orthonormal basis o

$$\tilde{\mathbf{x}} = \frac{\mathbf{b}_j^T \mathbf{x}}{\|\mathbf{b}_j\|^2} \mathbf{b}_j = \mathbf{b}_j \mathbf{b}_j^T \mathbf{x} \in \mathbb{R}^D$$

is the orthogonal projection of  $\mathbf{x}$  onto the subspace spanned by the  $j$ th basis vector, and  $z_j = \mathbf{b}_j^T \mathbf{x}$  is the coordinate of this projection with respect to the basis vector  $\mathbf{b}_j$  that spans that subspace since  $z_j \mathbf{b}_j = \tilde{\mathbf{x}}$ .

<https://eduassistpro.github.io/>

- More generally, if we aim to project onto a al subspace of  $\mathbb{R}^D$ , we obtain the orthogonal projection of  $\mathbf{x}$  onto t onal subspace with orthonormal basis vectors  $\mathbf{b}_1, \dots, \mathbf{b}_M$  as

$$\tilde{\mathbf{x}} = \mathbf{B} (\underbrace{\mathbf{B}^T \mathbf{B}}_{= \mathbf{I}})^{-1} \mathbf{B}^T \mathbf{x} = \mathbf{B} \mathbf{B}^T \mathbf{x}$$

where we defined  $\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$ . The coordinates of this projection with respect to the ordered basis  $(\mathbf{b}_1, \dots, \mathbf{b}_M)$  are  $\mathbf{z} := \mathbf{B}^T \mathbf{x}$

- Although  $\tilde{\mathbf{x}} \in \mathbb{R}^D$ , we only need  $M$  coordinates to represent  $\tilde{\mathbf{x}}$ . The other  $D - M$  coordinates with respect to the basis vectors  $(\mathbf{b}_{M+1}, \dots, \mathbf{b}_D)$  are always 0

## Assignment Project Exam Help

### 10.3.3 Finding the Basis of the Principal Subspace

- So far we have shown that for a given ONB optimal coordinates of  $\tilde{x}$  by an orthogonal projection onto the principal subspace following, we will determine what the **best basis** is.
- Recall the optimal coordinates of  $\tilde{x}$  given ONB is

$$z_{in} = \mathbf{x}_n^T \mathbf{b}_i = \mathbf{b}_i^T \mathbf{x}_n$$

- We have

## Assignment Project Exam Help

$$\mathbf{x}_n = \sum_{m=1}^M z_{mn} \mathbf{b}_m = \sum_{m=1}^M (\mathbf{x}_n^T \mathbf{b}_m) \mathbf{b}_m$$

- We now exploit the symmetry <https://eduassistpro.github.io/>

$$\tilde{\mathbf{x}}_n = \sum_{m=1}^M (\mathbf{b}_m^T \mathbf{x}_n) \mathbf{b}_m = \sum_{m=1}^M (\mathbf{x}_n^T \mathbf{b}_m) \mathbf{b}_m = \left( \sum_{m=1}^M \mathbf{b}_m \mathbf{b}_m^T \right) \mathbf{x}_n$$

- Since we can generally write the original data point  $\mathbf{x}_n$  as a linear combination of all basis vectors, it holds that

$$\begin{aligned} \mathbf{x}_n &= \sum_{d=1}^D z_{dn} \mathbf{b}_d = \sum_{d=1}^D (\mathbf{x}_n^T \mathbf{b}_d) \mathbf{b}_d = \left( \sum_{d=1}^D \mathbf{b}_d \mathbf{b}_d^T \right) \mathbf{x}_n \\ &= \left( \sum_{m=1}^M \mathbf{b}_m \mathbf{b}_m^T \right) \mathbf{x}_n + \left( \sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^T \right) \mathbf{x}_n \end{aligned}$$

where we split the sum with  $D$  terms into a sum over  $M$  and a sum over  $D - M$  terms.

# Assignment Project Exam Help

- With these results, the displacement vector  $\mathbf{x}_n - \tilde{\mathbf{x}}_n$ , i.e., the difference vector between the original data point and its projection, is

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \left( \sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^T \right) \mathbf{x}_n = \sum_{j=M+1}^D (\mathbf{x}_n^T \mathbf{b}_j) \mathbf{b}_j$$

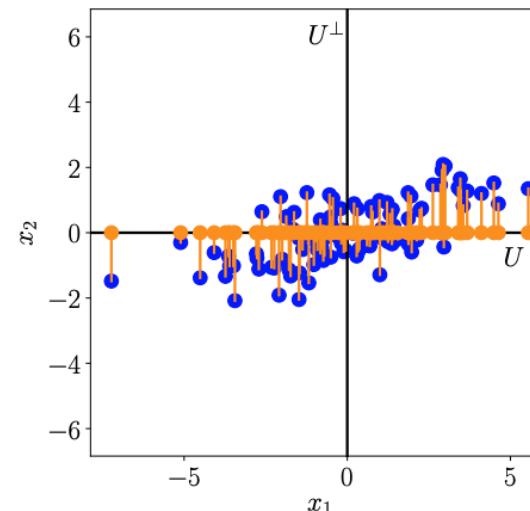
- The displacement vector  $\mathbf{x}_n - \tilde{\mathbf{x}}_n$  is exactly the projection of the data point onto the orthogonal complement of the principal subspace.

# Assignment Project Exam Help

- $\mathbf{x}_n - \tilde{\mathbf{x}}_n$  lies in the subspace  $U^\perp$  of the principal subspace.
- We identify the matrix  $\sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^T$  above as the projection matrix that performs this projection.

<https://eduassistpro.github.io/>

Orthogonal projection and displacement vectors. When projecting data points  $\mathbf{x}_n$  (blue) onto subspace  $U_1$ , we obtain  $\tilde{\mathbf{x}}_n$  (orange). The displacement vector  $\mathbf{x}_n - \tilde{\mathbf{x}}_n$  lies completely in the orthogonal complement  $U_2$  of  $U_1$ .



# Assignment Project Exam Help

- Now we reformulate the loss function.

$$J_M = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 + \frac{1}{N} \sum_{n=1}^N \left\| \sum_{j=M+1}^D (\mathbf{b}_j^\top x_n) \mathbf{b}_j \right\|^2$$

- We explicitly compute the squared norm and exploit the fact that the  $\mathbf{b}_j$  form an ONB:

$$J_M = \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D (\mathbf{b}_j^\top x_n)^2 + \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{b}_j x_n$$

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

where we exploited the symmetry of the dot product. The last step to write  $\mathbf{b}_j^\top x_n = x_n^\top \mathbf{b}_j$ . We now swap the sums and obtain

$$\begin{aligned} J_M &= \sum_{j=M+1}^D \mathbf{b}_j^\top \underbrace{\left( \frac{1}{N} \sum_{n=1}^N x_n x_n^\top \right)}_{=: S} \mathbf{b}_j = \sum_{j=M+1}^D \mathbf{b}_j^\top S \mathbf{b}_j \\ &= \sum_{j=M+1}^D \text{tr}(\mathbf{b}_j^\top S \mathbf{b}_j) = \sum_{j=M+1}^D \text{tr}(S \mathbf{b}_j \mathbf{b}_j^\top) = \text{tr} \left( \left( \sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^\top \right) S \right) \end{aligned}$$

where we exploited the property that the trace operator  $\text{tr}(\cdot)$  is linear and invariant to cyclic permutations of its arguments

# Assignment Project Exam Help

Add WeChat edu\_assist\_pro

$$J_M = \sum_{j=M+1}^D b_j^T S b_j = \left( \begin{pmatrix} & & \\ & \dots & \\ & & j b_j^T \\ & & \dots \\ & & & b_{M+1}^T \end{pmatrix} S \right)$$

projection matrix

- The loss is formulated as the covariance matrix of the data, projected onto the orthogonal complement of the principal subspace.
- Minimizing the average variance is therefore equivalent to minimizing the variance onto the subspace we ignore, i.e., the orthogonal subspace.
- Equivalently, we maximize the variance that we retain in the principal subspace, which links the projection immediately to the maximum-variance formulation of PCA in Section 10.2.
- In Section 10.2, the average squared reconstruction error, when projecting onto the  $M$ -dimensional principal subspace, is

$$J_M = \sum_{j=M+1}^D \lambda_j$$

- where  $\lambda_j$  are the eigenvalues of the data covariance matrix.

# Assignment Project Exam Help

Add WeChat `edu_assist_pro`

$$J_M = \sum_{j=M+1}^D \lambda_j$$

- To minimize it, we need to select the smallest  $D - M$  eigenvalues. Their corresponding eigenvectors are the basis of the orthogonal complement of the principal subspace.
- Consequently, this means the principal subspace comprises the eigenvectors  $b_1, \dots, b_M$  that are associated with the largest  $M$  eigenvalues of the data covariance matrix.

# Assignment Project Exam Help

## 10.5 PCA in High Dimensions

Add WeChat `edu_assist_pro`

- In order to do PCA, we need to compute the data covariance matrix  $\mathbf{S}$
- In  $D$  dimensions,  $\mathbf{S}$  is a  $D \times D$  matrix.
- Computing the eigenvalues and eigenvectors of this matrix is computationally expensive as it scales ~~cubically~~ in  $D$ .
- Therefore, PCA will be ~~solutions~~
- For example, if  $\mathbf{x}_n$  are in <https://eduassistpro.github.io/>, we would need to compute the eigendecomposition of a  $10,000 \times 10$ .
- We provide a solution to this problem for ~~it we have substantially~~ fewer data points than dimensions, i.e.,  $N \ll D$
- Assume we have a centered dataset  $\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x}_n \in \mathbb{R}^D$ . Then the data covariance matrix is given as

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^T \in \mathbb{R}^{D \times D}$$

where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$  is a  $D \times N$  matrix whose columns are the data points.

# Assignment Project Exam Help

- We now assume that  $N \ll D$ , i.e., the number of data points is smaller than the dimensionality of the data. [Add WeChat edu\\_assist\\_pro](#)
- With  $N \ll D$  data points, the rank of the covariance matrix  $S$  is at most  $N$ , so it has at least  $D - N$  eigenvalues that are 0.
- Intuitively, this means that there are some redundancies. In the following, we will exploit this and turn the  $D \times D$  covariance matrix into an  $N \times N$  covariance matrix whose eigenvalues are all positive. [Assignment Project Exam Help](#)
- In PCA, we ended up

<https://eduassistpro.github.io/>

where  $b_m$  is a basis vector of the principal component space. Let us rewrite this equation a bit: With  $S = \frac{1}{N}XX^T \in \mathbb{R}^{D \times D}$ , we

$$Sb_m = \frac{1}{N}XX^Tb_m = \lambda_m b_m$$

- We now multiply  $X^T \in \mathbb{R}^{N \times D}$  from the left-hand side, which yields

$$\underbrace{\frac{1}{N}X^TX}_{N \times N} \underbrace{X^Tb_m}_{=: c_m} = \lambda_m X^Tb_m \Leftrightarrow \frac{1}{N}X^TXc_m = \lambda_m c_m$$

# Assignment Project Exam Help

Add WeChat  $\frac{1}{N} \mathbf{X}^T \mathbf{X}$  [edu\\_assist\\_pro](#)

- We get a new eigenvector/eigenvalue equation:  $\lambda_m$  remains eigenvalue, which confirms our results from exercise 4.11 that the nonzero eigenvalues of  $\mathbf{X} \mathbf{X}^T$  equal the nonzero eigenvalues of  $\mathbf{X}^T \mathbf{X}$ .
- We obtain the eigenvector of the matrix  $\frac{1}{N} \mathbf{X}^T \mathbf{X} \in \mathbb{R}^{N \times N}$  associated with  $\lambda_m$  as  $\mathbf{c}_m := \mathbf{X}^T \mathbf{b}_m$ .

# Assignment Project Exam Help

- This also implies that  $\frac{1}{N} \mathbf{X}^T \mathbf{X}$  has the same (nonzero) eigenvalues as the data covariance matrix  $\mathbf{S}$ .
- But  $\mathbf{X}^T \mathbf{X}$  now an  $N \times N$  matrix, so we can recover the eigenvalues and eigenvectors much more efficiently than for the  $\mathbf{S}$  data covariance matrix.
- Now that we have the eigenvectors of  $\frac{1}{N} \mathbf{X}^T \mathbf{X}$ , we are going to recover the original eigenvectors, which we still need for PCA. Currently, we know the eigenvectors of  $\frac{1}{N} \mathbf{X}^T \mathbf{X}$ . If we left-multiply our eigenvalue/ eigenvector equation with  $\mathbf{X}$ , we get

$$\frac{1}{N} \mathbf{X} \mathbf{X}^T \mathbf{X} \mathbf{c}_m = \lambda_m \mathbf{X} \mathbf{c}_m$$

$\underbrace{\phantom{\mathbf{X} \mathbf{X}^T \mathbf{X}}}_{\mathbf{S}}$

and we recover the data covariance matrix again. This now also means that we recover  $\mathbf{X} \mathbf{c}_m$  as an eigenvector of  $\mathbf{S}$ .

# 10.6 Key Steps of PCA in Practice

Assignment Project Exam Help  
Add WeChat edu\_assist\_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

eigendecomposition

# Assignment Project Exam Help

- **Step 1. Mean subtraction** Divide the data points by the standard deviation  $\sigma_d$  of the dataset for every dimension  $d = 1, \dots, D$ . Now the data has variance 1 along each axis.
- We center the data by computing the mean  $\mu$  of the dataset and subtracting it from every single data point. This ensures that the dataset has mean 0.

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

- Step 3. Eigendecomposition of the  $\mathbf{X}^T \mathbf{X}$  matrix  
Add WeChat `edu_assist_pro`
- Compute the data covariance matrix values and corresponding eigenvectors. The longer vector (larger eigenvalue) spans the principal subspace  $U$

# Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

# Assignment Project Exam Help

- 4. **Projection** We can project any  $\mathbf{x} \in \mathbb{R}^D$  onto the principal subspace: To get this right, we need  $\mathbf{x}_*$  using the mean  $\mu_d$  and standard deviation  $\sigma_d$  of the training data in the  $d$ th dimension, respectively, so that

$$x_*^{(d)} \leftarrow \frac{x_*^{(d)} - \mu_d}{\sigma_d}, \quad d = 1, \dots, D$$

Assignment Project Exam Help

where  $x_*^{(d)}$  is the  $d$ th c

- We obtain the projection <https://eduassistpro.github.io/>

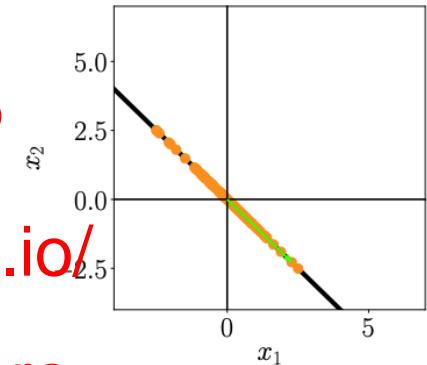
with coordinates

$\tilde{\mathbf{x}}_* = \mathbf{B}\mathbf{B}^\top \mathbf{x}_*$

$$\mathbf{z}_* = \mathbf{B}^\top \mathbf{x}_*$$

with respect to the basis of the principal subspace. Here,  $\mathbf{B}$  is the matrix that contains the eigenvectors that are associated with the largest eigenvalues of the data covariance matrix as columns.

- Note that PCA returns the coordinates  $\mathbf{z}_*$ , not the projections of  $\mathbf{x}_*$ .



(e) Step 4: Project data onto the principal subspace.

# Assignment Project Exam Help

- Having standardized our dataset  $\tilde{x}_*$  yields the projections in the context of the **standardized datas** **Add WeChat edu\_assist\_pro**
- To obtain our projection in the original data space (i.e., before standardization), we need to undo the standardization: multiply by the standard deviation before adding the mean.
- We obtain **Assignment Project Exam Help**  
 $\tilde{x}_*^{(d)}$        $(d)$        $\dots, D$
- Figure 10.10(f) illustrates **https://eduassistpro.github.io/** all data space.

**Add WeChat edu\_assist\_pro**