

Assignment Project Exam Help

Add WeChat powcoder

# Dimensionality Reduction with Principal Component Analysis

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Liang Zheng

Australian National University

liang.zheng@anu.edu.au

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

Add WeChat powcoder

Assignment Project Exam Help

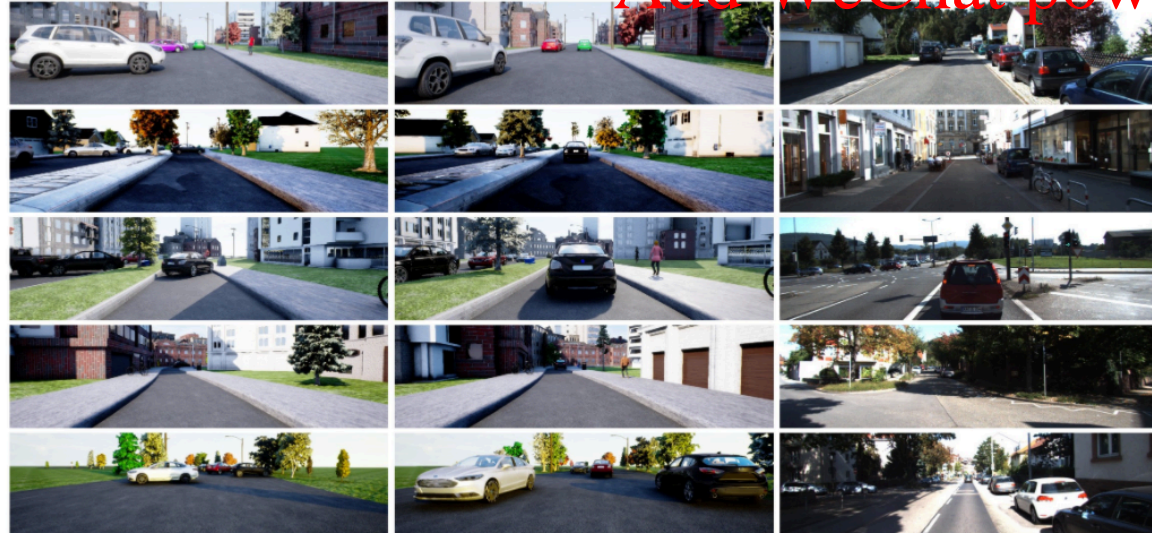
<https://powcoder.com>

Input Prob. Grammar

Meta-Sim

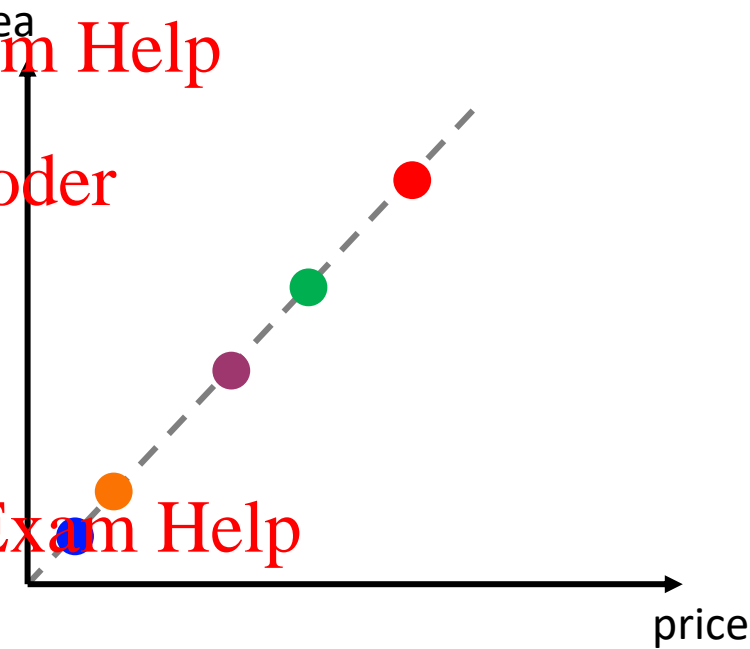
KITTI Dataset

Add WeChat powcoder



# Idea of PCA

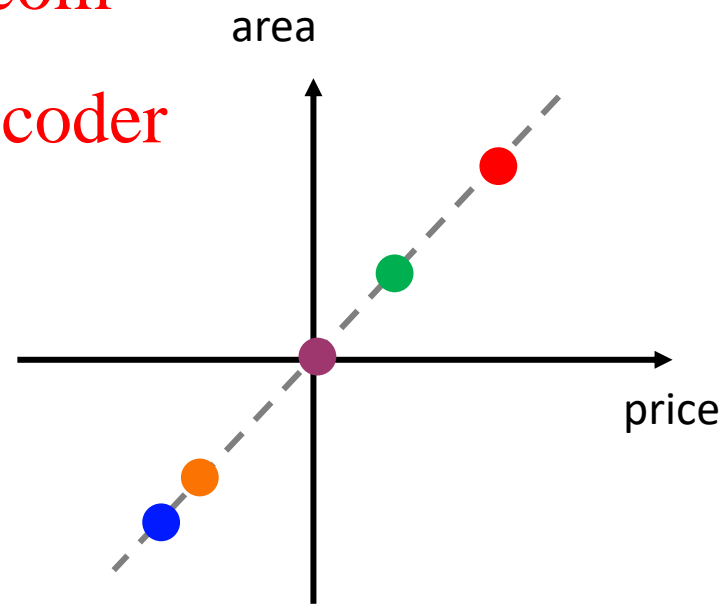
	House price (million)	House area (100m <sup>2</sup> )
a	10	10
b	2	2
c	7	7
d	1	1
e	5	5



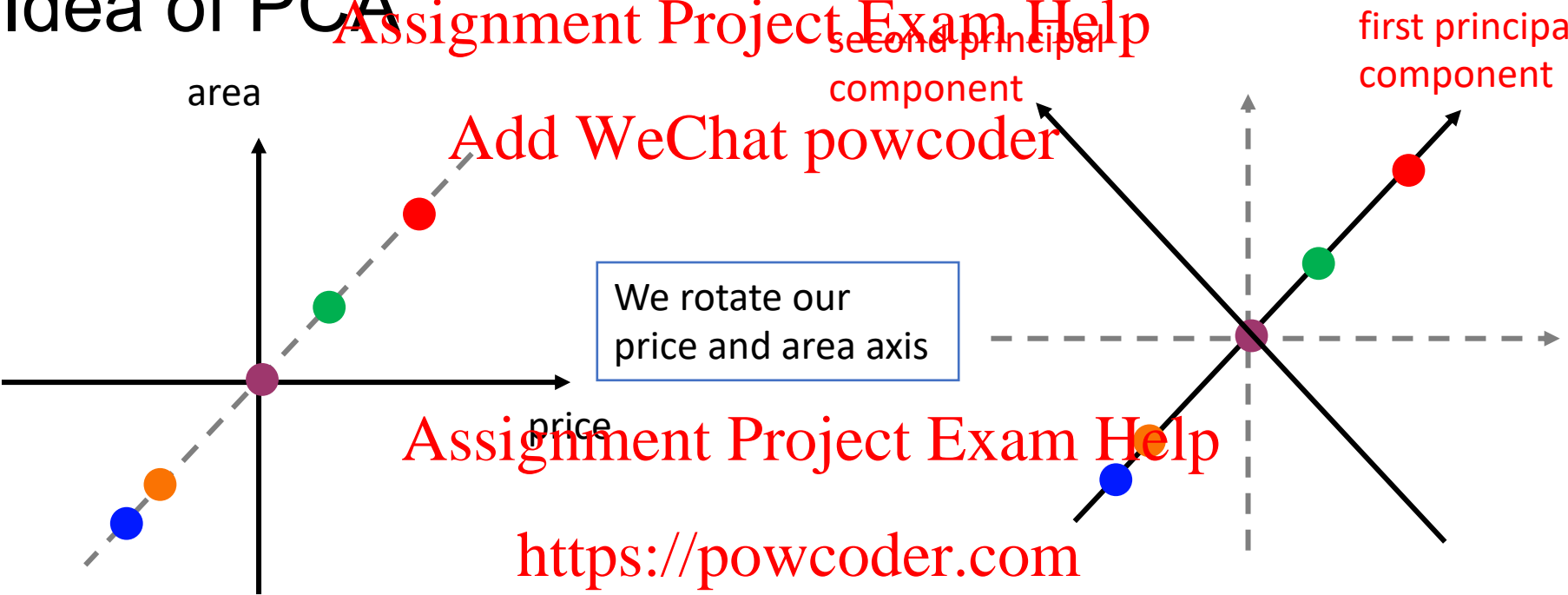
<https://powcoder.com>

We subtract means from data points

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



# Idea of PCA



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

	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

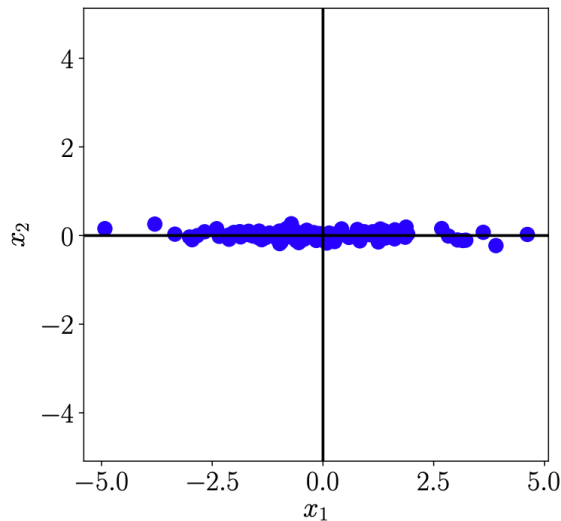
## Motivation

Add WeChat powcoder

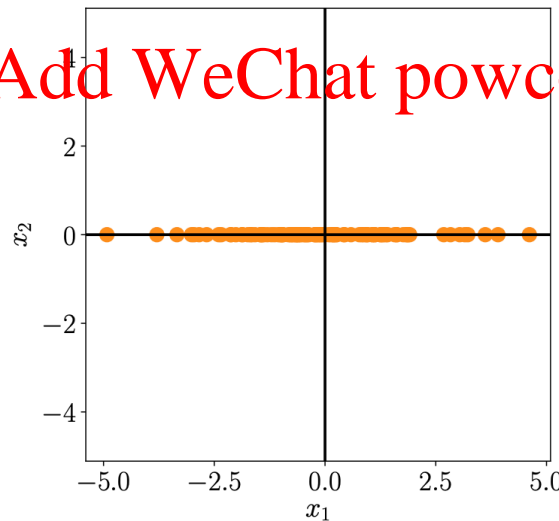
- High-dimensional data, such as images, is hard to analyze, interpretate, and visualize, and expensive to store.
- Good news
- high-dimensional data is often overcomplete, i.e., many dimensions are redundant and can be explained by a combination of other dimensions
- Furthermore, dimensions in high-dimensional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

Assignment Project Exam Help

<https://powcoder.com>



(a) Dataset with  $x_1$  and  $x_2$  coordinates.



(b) Compressed dataset where only the  $x_1$  coordinate is relevant.

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

Add WeChat powcoder

- In PCA, we are interested in finding projections  $\tilde{x}_n$  of data points  $x_n$  that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality
- We consider an i.i.d. dataset  $\mathcal{X} = \{x_1, \dots, x_N\}$ ,  $x_n \in \mathbb{R}^D$ , with mean  $\mathbf{0}$  that possesses the data covariance matrix

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

Assignment Project Exam Help

<https://powcoder.com>

- We assume there exists a low-dimensional compressed representation (code)

$$z_n = B^T x_n \in \mathbb{R}^M$$

of  $x_n$ , where we define the projection matrix

$$B := [b_1, \dots, b_M] \in \mathbb{R}^{D \times M}$$

Add WeChat powcoder

# Assignment Project Exam Help

Add WeChat powcoder

- **Example (Coordinate Representation/Code)**
- Consider  $\mathbb{R}^2$  with the canonical basis  $\mathbf{e}_1 = [1, 0]^T$ ,  $\mathbf{e}_2 = [0, 1]^T$ .
- $\mathbf{x} \in \mathbb{R}^2$  can be represented as a linear combination of these basis vectors, e.g.,

$$\begin{bmatrix} 5 \\ 3 \end{bmatrix} = 5\mathbf{e}_1 + 3\mathbf{e}_2$$

Assignment Project Exam Help

- However, when we consider vectors of the form

$$\tilde{\mathbf{x}} = \begin{bmatrix} 0 \\ z \end{bmatrix} \in \mathbb{R}^2, \quad z \in \mathbb{R}$$

<https://powcoder.com>

they can always be written as  $0\mathbf{e}_1 + z\mathbf{e}_2$ .

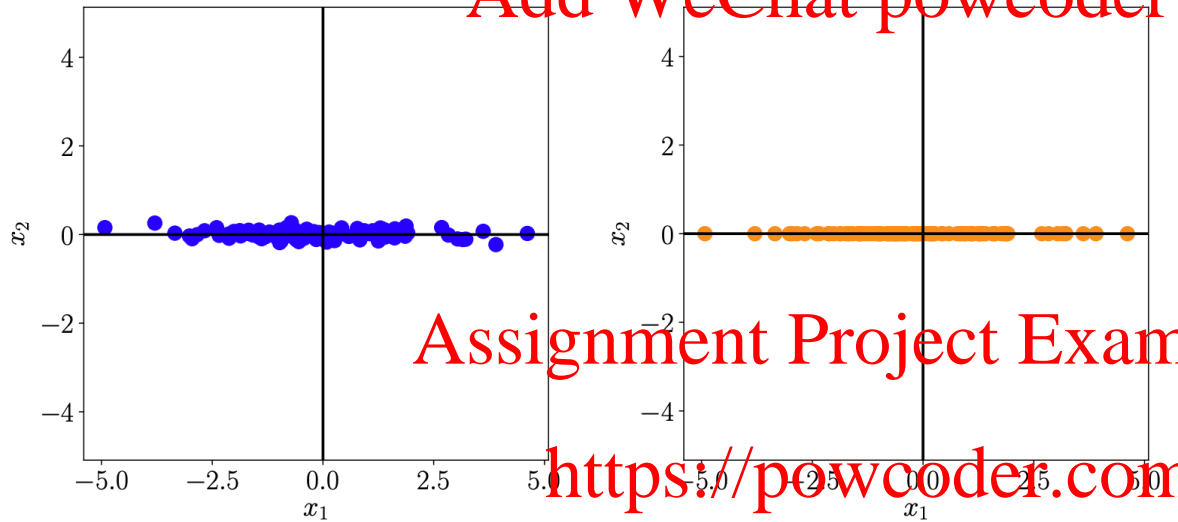
Add WeChat powcoder

- To represent these vectors it is sufficient to store the coordinate/code  $z$  of  $\tilde{\mathbf{x}}$  with respect to the  $\mathbf{e}_2$  vector.



# 10.2 PCA from Maximum Variance Perspective

Add WeChat powcoder



Assignment Project Exam Help

<https://powcoder.com>

(a) Dataset with  $x_1$  and  $x_2$  coordinates.

(b) Compressed dataset where only the  $x_1$  coordinate is relevant.

Add WeChat powcoder

- 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 powcoder

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

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

Assignment Project Exam Help

- Let us assume that  $\mu$  is the mean of the data. Using the properties of the variance, we obtain

<https://powcoder.com>

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

Add WeChat powcoder

- That is, the variance of the low-dimensional code does not depend on the mean of the data.
- With this assumption the mean of the low-dimensional code is also  $\mathbf{0}$  since

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

# Assignment Project Exam Help

- To maximize the variance of the low-dimensional code, we first seek a single vector  $\mathbf{b}_1 \in \mathbb{R}^D$  that maximizes the variance of the projected data, i.e., we aim to maximize the variance of the first coordinate  $z_1$  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,

$$z_{1n} = \mathbf{b}_1^T \mathbf{x}_n$$

i.e., it is the coordinate of the orthogonal projection of  $\mathbf{x}_n$  onto the one-dimensional subspace spanned by  $\mathbf{b}_1$ . We substitute  $z_{1n}$  into  $V_1$  and obtain,

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

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

$$\begin{aligned} & \max_{\mathbf{b}_1} \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1 \\ & \text{subject to } \|\mathbf{b}_1\|^2 = 1 \end{aligned}$$

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

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

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

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

- Setting these partial derivatives to 0 gives us the relations

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

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

- 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**.

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

Add WeChat powcoder

- Assume we have found the  $m-1$  eigenvectors of  $S$  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,

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

https://powcoder.com

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

Add WeChat powcoder

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

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

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

# Assignment Project Exam Help

Add WeChat powcoder

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

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

- $\mathbf{b}_m$  is not only an eigenvector of  $\mathbf{S}$  but also of  $\hat{\mathbf{S}}$ .
- Specifically,  $\lambda_m$  is the largest eigenvalue of  $\mathbf{S}$  and  $\lambda_m$  is the  $m$ th largest eigenvalue of  $\hat{\mathbf{S}}$ , and both have the associated eigenvector  $\mathbf{b}_m$ .
- Moreover,  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$  are also eigenvectors of  $\mathbf{S}$ , but they 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.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Assignment Project Exam Help

## MNIST dataset

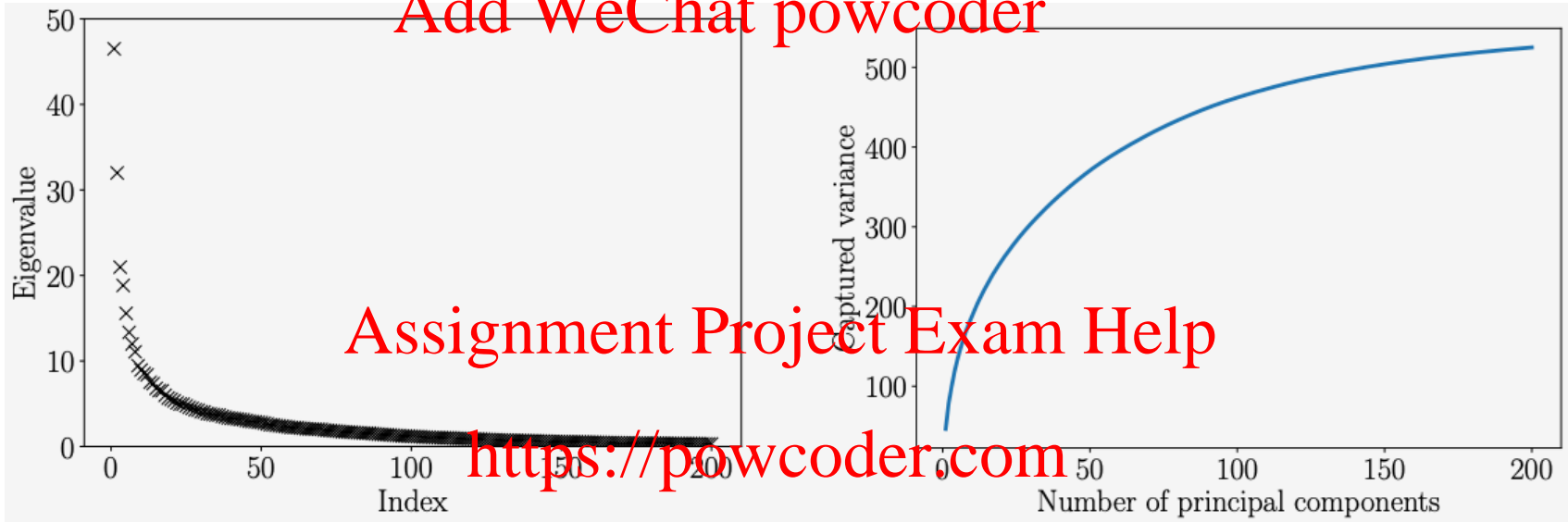
Add WeChat powcoder

- 60,000 examples of handwritten digits 0 through 9.
- 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  $x \in \mathbb{R}^{784}$



# Example - Eigenvalues of MNIST digit “8”

Add WeChat powcoder



Assignment Project Exam Help

<https://powcoder.com>

(a) Top 200 largest eigenvalues

(b) Variance captured by the principal components.

Add WeChat powcoder

- 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 powcoder

- To find an  $M$ -dimensional subspace of  $\mathbb{R}^D$  that retains as much information as possible,
- We choose the columns of  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$  as the  $M$  eigenvectors of the data covariance matrix  $\mathbf{S}$  that are associated with the  $M$  largest eigenvalues.

## Assignment Project Exam Help

- The maximum amount of variance PCA can capture with the first  $M$  principal components is

<https://powcoder.com>

$V_M = \sum_{m=1}^M \lambda_m$   
Add WeChat powcoder

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

- 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 powcoder

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

$$\max_{b_1} b_1^T S b_1$$

$$\text{subject to } \|b_1\|^2 = 1$$

Assignment Project Exam Help

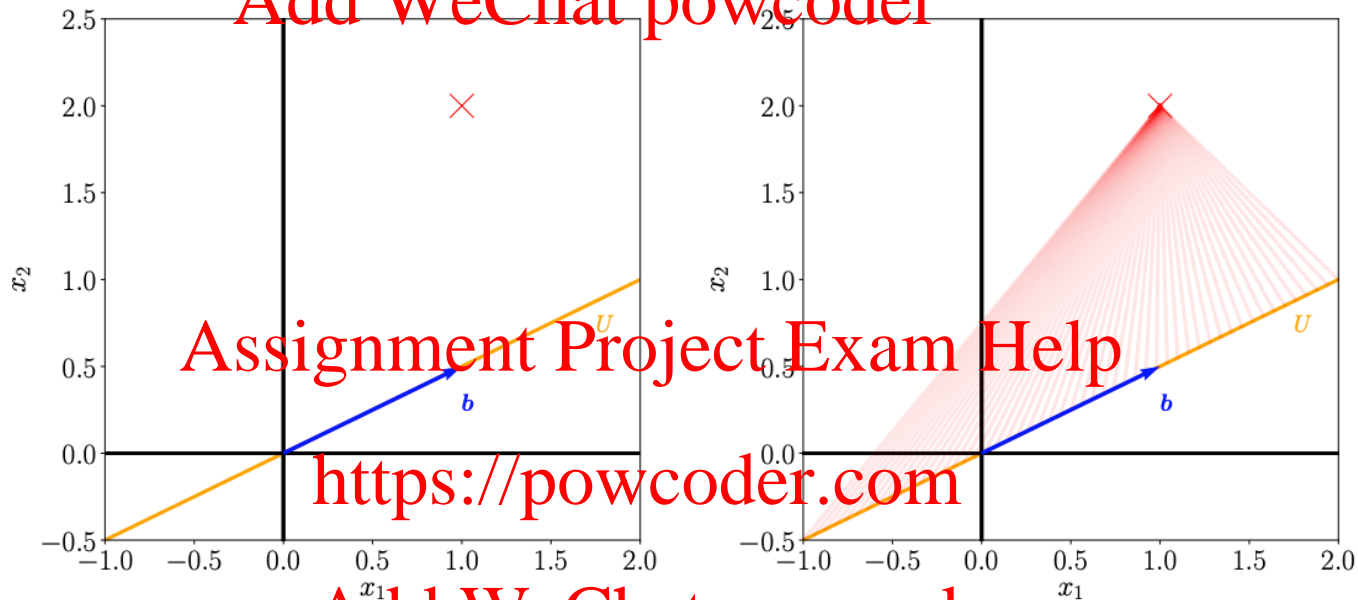
- Alternatively, we derive PCA as an algorithm that directly minimizes the average reconstruction error

<https://powcoder.com>

Add WeChat powcoder

# 10.3.1 Setting and Objective

Add WeChat powcoder



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

(a) A vector  $x \in \mathbb{R}^2$  (red cross) shall be projected onto a one-dimensional subspace  $U \subseteq \mathbb{R}^2$  spanned by  $b$

(b) Differences  $x - \tilde{x}_i$  for 50 different  $\tilde{x}_i$  are shown by the red lines

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

# Assignment Project Exam Help

## Add WeChat powcoder

- 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  $\zeta_d \in \mathbb{R}$ .

<https://powcoder.com>

- We aim to find vectors  $\tilde{\mathbf{x}} \in \mathbb{R}^D$ , which live in an intrinsically lower-dimensional subspace  $U \subseteq \mathbb{R}^D$ ,  $\dim(U) = M$ , so that

## Add WeChat powcoder

$$\tilde{\mathbf{x}} = \sum_{m=1}^M z_m \mathbf{b}_m \in U \subseteq \mathbb{R}^D$$

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

# Assignment Project Exam Help

Add WeChat powcoder

# Assignment Project Exam Help

$$\tilde{\mathbf{x}}_n := \sum_{m=1}^M z_{mn} \mathbf{b}_m = \mathbf{B} \mathbf{z}_n \in \mathbb{R}^D$$

<https://powcoder.com>

Add WeChat powcoder

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

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

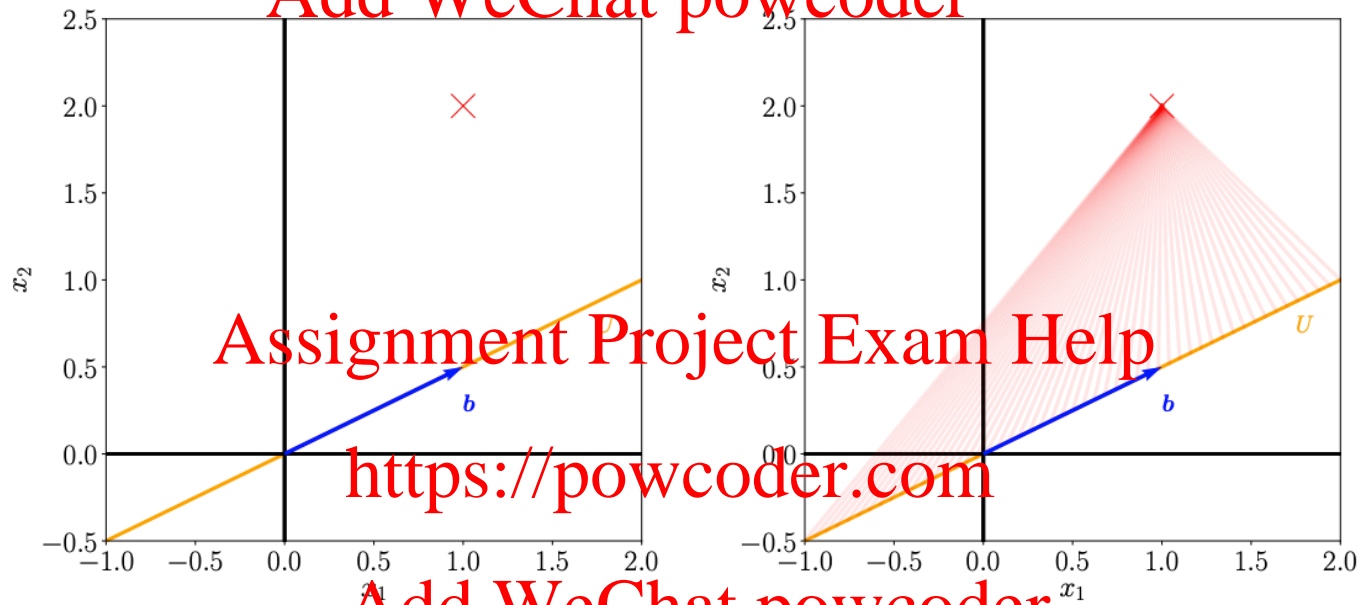
- We want to have  $\tilde{\mathbf{x}}_n$  as similar to  $\mathbf{x}_n$  as possible.
- We define our objective as minimizing the average squared Euclidean distance (**reconstruction error**)

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

- 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.

## 10.3.2 Finding Optimal Coordinates

Add WeChat powcoder



Add WeChat powcoder

(a) A vector  $x \in \mathbb{R}^2$  (red cross) shall be projected onto a one-dimensional subspace  $U \subseteq \mathbb{R}^2$  spanned by  $b$

(b) Differences  $x - \tilde{x}_i$  for 50 different  $\tilde{x}_i$  are shown by the red lines

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

# Assignment Project Exam Help

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

Add WeChat powcoder

- Given an ONB  $(\mathbf{b}_1, \dots, \mathbf{b}_M)$  of  $U \subseteq \mathbb{R}^D$ , to find the optimal coordinates  $\mathbf{z}_m$  with respect to this basis, we calculate the partial derivatives

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

Assignment Project Exam Help

$$\tilde{\mathbf{x}}_n := \sum_{m=1}^M z_{mn} \mathbf{b}_m = \mathbf{B} \mathbf{z}_n \in \mathbb{R}^D$$

$$\frac{\partial \tilde{\mathbf{x}}_n}{\partial z_{in}} = \frac{\partial}{\partial z_{in}} \left( \sum_{m=1}^M z_{mn} \mathbf{b}_m \right) = \mathbf{b}_i$$

Add WeChat powcoder

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

$$\begin{aligned} \frac{\partial J_M}{\partial z_{in}} &= -\frac{2}{N} (\mathbf{x}_n - \tilde{\mathbf{x}}_n)^T \mathbf{b}_i = -\frac{2}{N} \left( \mathbf{x}_n - \sum_{m=1}^M z_{mn} \mathbf{b}_m \right)^T \mathbf{b}_i \\ &\stackrel{\text{ONB}}{=} -\frac{2}{N} (\mathbf{x}_n^T \mathbf{b}_i - z_{in} \mathbf{b}_i^T \mathbf{b}_i) = -\frac{2}{N} (\mathbf{x}_n^T \mathbf{b}_i - z_{in}) \end{aligned}$$



# Assignment Project Exam Help

Add WeChat powcoder

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

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

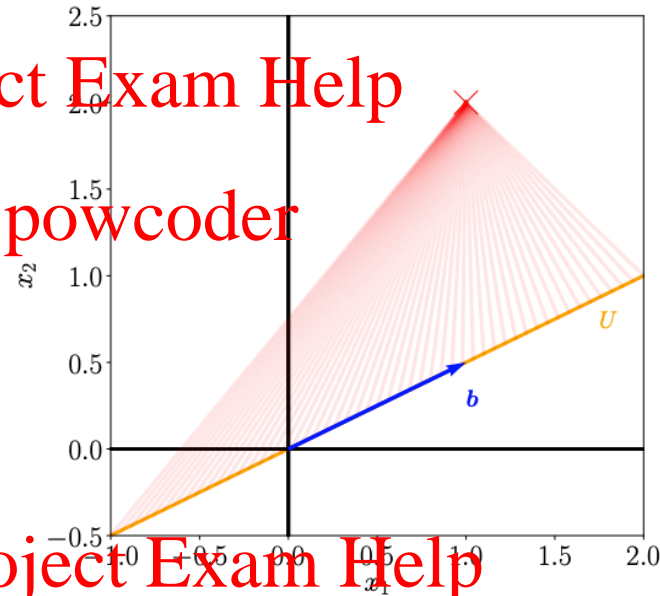
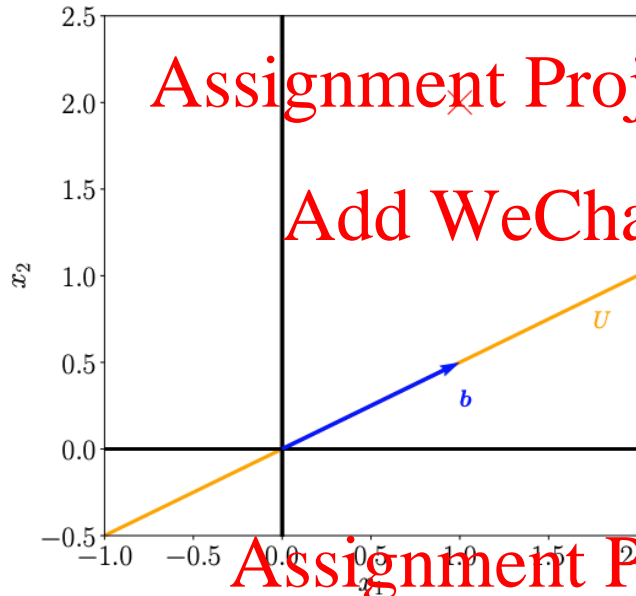
for  $i = 1, \dots, M$ , and  $n = 1, \dots, N$ .

Assignment Project Exam Help

- The optimal coordinates  $z_{in}$  of the projection  $\tilde{\mathbf{x}}_n$  are the coordinates of the orthogonal projection of the original data point  $\mathbf{x}_n$  onto the one-dimensional subspace that is spanned by  $\mathbf{b}_i$ .
- The optimal linear projection  $\tilde{\mathbf{x}}_n$  of  $\mathbf{x}_n$  is an orthogonal 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.

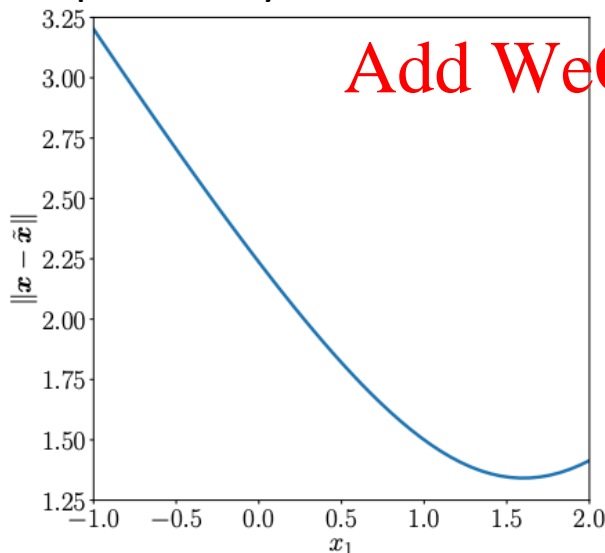
https://powcoder.com

Add WeChat powcoder

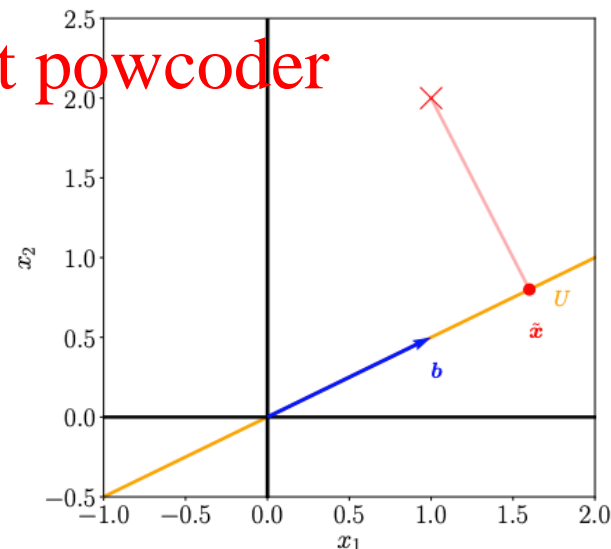


(a) A vector  $x \in \mathbb{R}^2$  (red cross) shall be projected onto a one-dimensional subspace  $U \subseteq \mathbb{R}^2$  spanned by  $b$

(b) Differences  $x - \tilde{x}_i$  for 50 different  $\tilde{x}_i$  are shown by the red lines



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



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

# Assignment Project Exam Help

Add WeChat powcoder

- We briefly recap orthogonal projections from Section 3.8 (Analytic geometry).
- If  $(\mathbf{b}_1, \dots, \mathbf{b}_D)$  is an orthonormal basis of  $\mathbb{R}^D$  then

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

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

Assignment Project Exam Help

<https://powcoder.com>

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

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

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  $\mathbf{z} := \mathbf{B}^\top \mathbf{x}$

- 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

## 10.3.3 Finding the Basis of the Principal Subspace

Add WeChat powcoder

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

- Recall the optimal coordinates of  $\tilde{\mathbf{x}}$  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

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

<https://powcoder.com>

- We now exploit the symmetry of the dot product, which yields

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

- 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.

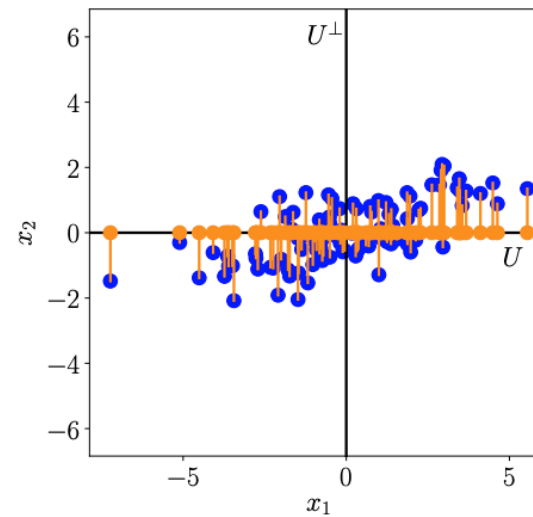
# Assignment Project Exam Help

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

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.

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

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

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

$$= \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D \mathbf{b}_j \mathbf{x}_n \mathbf{x}_n^T \mathbf{b}_j$$

<https://www.powcoder.com>

Add WeChat powcoder

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

$$J_M = \sum_{j=M+1}^D \mathbf{b}_j^T \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{b}_j = \sum_{j=M+1}^D \mathbf{b}_j^T \mathbf{S} \mathbf{b}_j$$

$$= \sum_{j=M+1}^D \text{tr}(\mathbf{b}_j^T \mathbf{S} \mathbf{b}_j) = \sum_{j=M+1}^D \text{tr}(\mathbf{S} \mathbf{b}_j \mathbf{b}_j^T) = \text{tr} \left( \underbrace{\left( \sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^T \right)}_{\text{projection matrix}} \mathbf{S} \right)$$

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

$$J_M = \sum_{j=M+1}^D \mathbf{b}_j^T \mathbf{S} \mathbf{b}_j = \text{tr} \left( \underbrace{\left( \sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^T \right)}_{\text{projection matrix}} \mathbf{S} \right)$$

- The loss is formulated as the covariance matrix of the data, projected onto the orthogonal complement of the principal subspace.
- Minimizing the average squared reconstruction error is therefore equivalent to minimizing the variance of the data when projected onto the subspace we ignore, i.e., the orthogonal complement of the principal subspace.
- Equivalently, we maximize the variance of the projection that we retain in the principal subspace, which links the projection loss 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 powcoder

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

- 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 that the basis of the principal subspace comprises the eigenvectors  $\mathbf{b}_1, \dots, \mathbf{b}_M$  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 powcoder

- 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 infeasible in very high dimensions
- For example, if  $\mathbf{x}_n$  are images with 10,000 pixels, we would need to compute the eigendecomposition of a  $10,000 \times 10,000$  matrix.
- We provide a solution to this problem for the case that 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

<https://powcoder.com>

Add WeChat powcoder

# 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.
- 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 with the eigenvector equation

$$S \mathbf{b}_m = \lambda_m \mathbf{b}_m, \quad m = 1, \dots, M$$

where  $\mathbf{b}_m$  is a basis vector of the principal subspace. Let us rewrite this equation a bit: With  $S = \frac{1}{N} \mathbf{X} \mathbf{X}^T \in \mathbb{R}^{D \times D}$ , we obtain

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

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

$$\frac{1}{N} \underbrace{\mathbf{X}^T \mathbf{X}}_{N \times N} \underbrace{\mathbf{X}^T \mathbf{b}_m}_{=: \mathbf{c}_m} = \lambda_m \mathbf{X}^T \mathbf{b}_m \Leftrightarrow \frac{1}{N} \mathbf{X}^T \mathbf{X} \mathbf{c}_m = \lambda_m \mathbf{c}_m$$

# Assignment Project Exam Help

Add WeChat powcoder

- 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 that we can compute the eigenvalues and eigenvectors much more efficiently than for the original  $D \times D$  data covariance matrix.

<https://powcoder.com>

Add WeChat powcoder

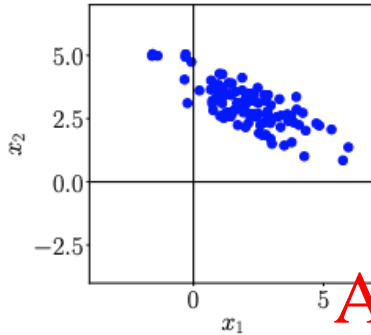
- Now that we have the eigenvectors of  $\frac{1}{N}\mathbf{X}^T\mathbf{X}$ , we are going to recover the original eigenvectors, which we still need for PCA. Currently, we know the eigenvectors of  $\frac{1}{N}\mathbf{X}^T\mathbf{X}$ . If we left-multiply our eigenvalue/ eigenvector equation with  $\mathbf{X}$ , we get

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

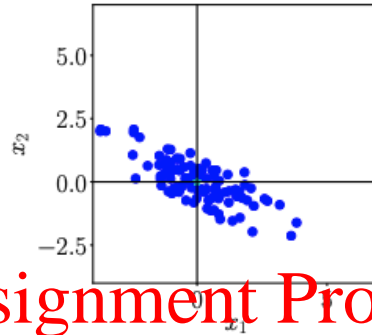
and we recover the data covariance matrix again. This now also means that we recover  $\mathbf{X}\mathbf{c}_m$  as an eigenvector of  $\mathbf{S}$ .

# 10.6 Key Steps of PCA in Practice

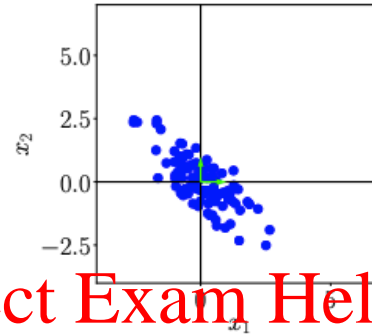
Add WeChat powcoder



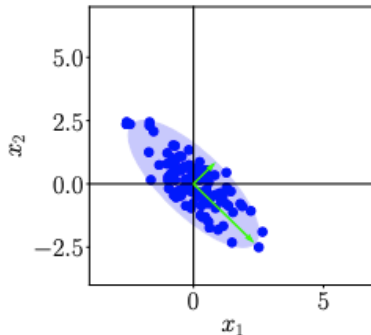
(a) Original dataset.



(b) Step 1: Centering by subtracting the mean from each data point.

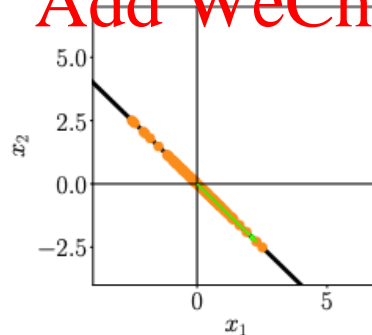


(c) Step 2: Dividing by the standard deviation to make the data unit free. Data has variance 1 along each axis.

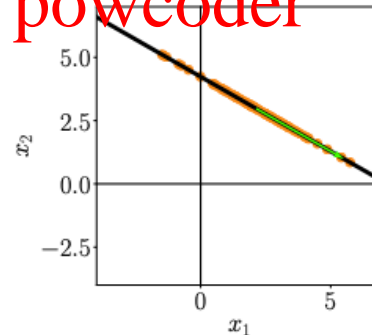


(d) Step 3: Compute eigenvalues and eigenvectors (arrows) of the data covariance matrix (ellipse).

eigendecomposition



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



(f) Undo the standardization and move projected data back into the original data space from (a).

<https://powcoder.com>

Add WeChat powcoder

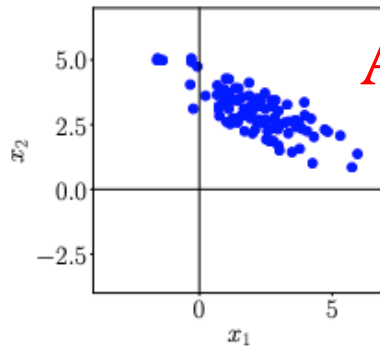
# Assignment Project Exam Help

Add WeChat powcoder

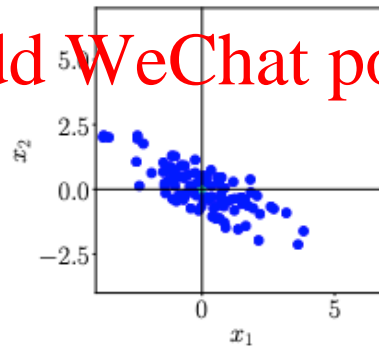
- **Step 1. Mean subtraction**
- We center the data by computing the mean  $\mu$  of the dataset and subtracting it from every single data point. This ensures that the dataset has mean 0.
- **Step 2. Standardization** Divide the data points by the standard deviation  $\sigma$  of the dataset for every dimension  $d = 1, \dots, D$ . Now the data has variance 1 along each axis.

<https://powcoder.com>

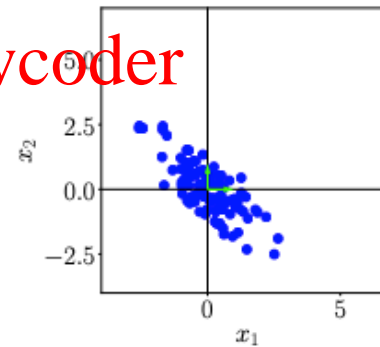
Add WeChat powcoder



(a) Original dataset.



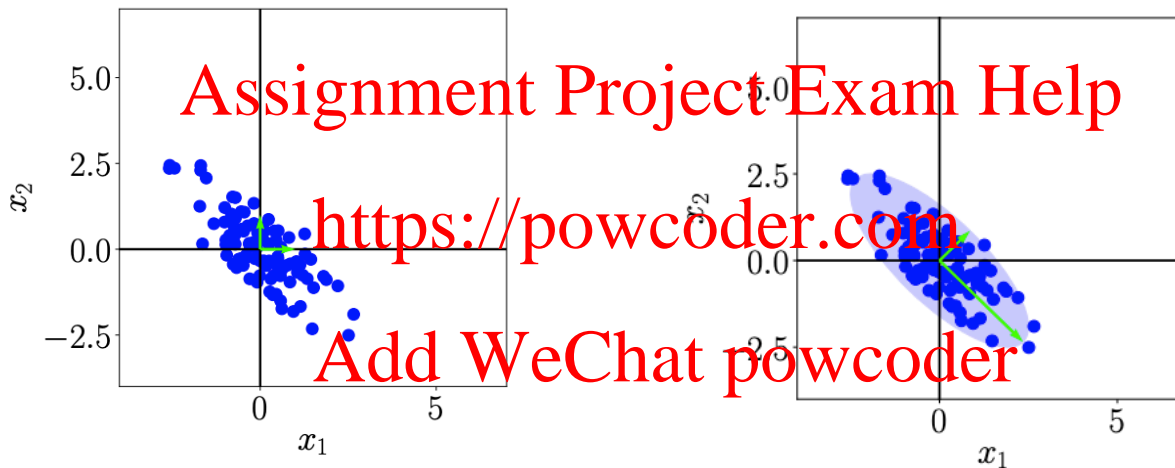
(b) Step 1: Centering by subtracting the mean from each data point.



(c) Step 2: Dividing by the standard deviation to make the data unit free. Data has variance 1 along each axis.

# Assignment Project Exam Help

- **Step 3. Eigendecomposition of the covariance matrix**
- Compute the data covariance matrix and its eigenvalues and corresponding eigenvectors. The longer vector (larger eigenvalue) spans the principal subspace  $U$



(c) Step 2: Dividing by the standard deviation to make the data unit free. Data has variance 1 along each axis.

(d) Step 3: Compute eigenvalues and eigenvectors (arrows) of the data covariance matrix (ellipse).



# Assignment Project Exam Help

Add WeChat powcoder

- **4. Projection** We can project any data point  $\mathbf{x}_* \in \mathbb{R}^D$  onto the principal subspace: To get this right, we need to standardize  $\mathbf{x}_*$  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 component of  $\mathbf{x}_*$ .

- We obtain the projection as <https://powcoder.com>

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

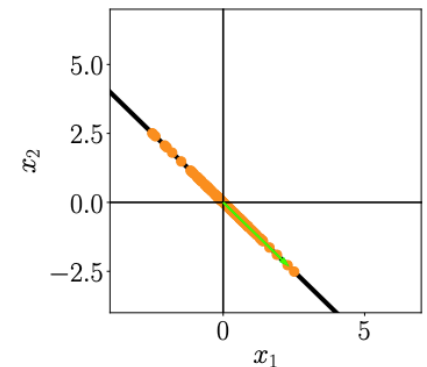
Add WeChat powcoder

with coordinates

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

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.

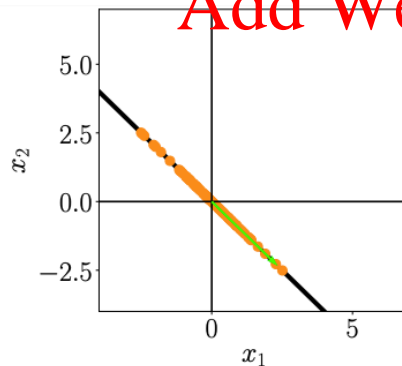
# Assignment Project Exam Help

- Having standardized our dataset,  $\tilde{\mathbf{x}}_* = \mathbf{B}\mathbf{B}^T \mathbf{x}_*$  only yields the projections in the context of the standardized dataset.
- 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.

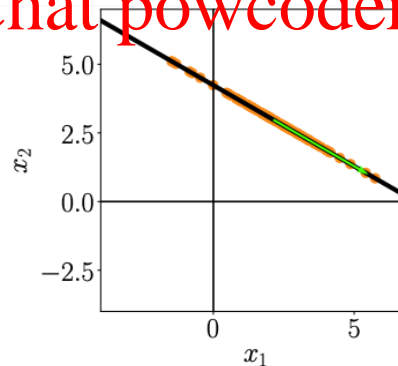
- We obtain

$$\tilde{\mathbf{x}}_*^{(d)} \leftarrow \tilde{\mathbf{x}}_*^{(d)} \sigma_d + \mu_d, \quad d = 1, \dots, D$$

- Figure 10.10(f) illustrates the projection in the original data space.



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



(f) Undo the standardization and move projected data back into the original data space from (a).