

Dimensionality Reduction with Principal Component Analysis

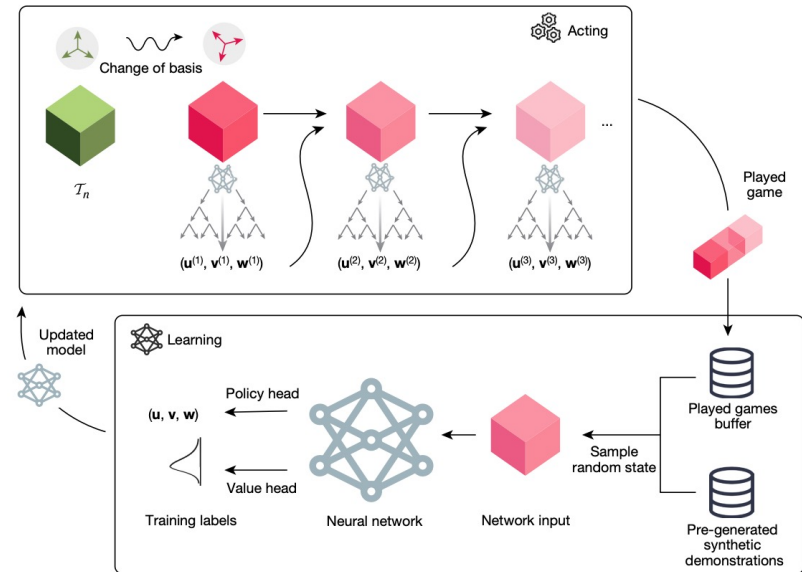
Liang Zheng
Australian National University
liang.zheng@anu.edu.au

Discovering faster matrix multiplication algorithms with reinforcement learning. Fawzi et al., Nature 2022

$$\begin{bmatrix} a_{1,1} & a_{1,2} \\ a_{2,1} & a_{2,2} \end{bmatrix} \times \begin{bmatrix} b_{1,1} & b_{1,2} \\ b_{2,1} & b_{2,2} \end{bmatrix} = \begin{bmatrix} c_{1,1} & c_{1,2} \\ c_{2,1} & c_{2,2} \end{bmatrix}$$

$h_1 = a_{1,1} b_{1,1}$	$h_1 = (a_{1,1} + a_{2,2})(b_{1,1} + b_{2,2})$
$h_2 = a_{1,1} b_{1,2}$	$h_2 = (a_{2,1} + a_{2,2}) b_{1,1}$
$h_3 = a_{1,2} b_{2,1}$	$h_3 = a_{1,1} (b_{1,2} - b_{2,2})$
$h_4 = a_{1,2} b_{2,2}$	$h_4 = a_{2,2} (-b_{1,1} + b_{2,1})$
$h_5 = a_{2,1} b_{1,1}$	$h_5 = (a_{1,1} + a_{1,2}) b_{2,2}$
$h_6 = a_{2,1} b_{1,2}$	$h_6 = (-a_{1,1} + a_{2,1})(b_{1,1} + b_{1,2})$
$h_7 = a_{2,2} b_{2,1}$	$h_7 = (a_{1,2} - a_{2,2})(b_{2,1} + b_{2,2})$
$h_8 = a_{2,2} b_{2,2}$	
$c_{1,1} = h_1 + h_3$	$c_{1,1} = h_1 + h_4 - h_5 + h_7$
$c_{1,2} = h_2 + h_4$	$c_{1,2} = h_3 + h_5$
$c_{2,1} = h_5 + h_7$	$c_{2,1} = h_2 + h_4$
$c_{2,2} = h_6 + h_8$	$c_{2,2} = h_1 - h_2 + h_3 + h_6$

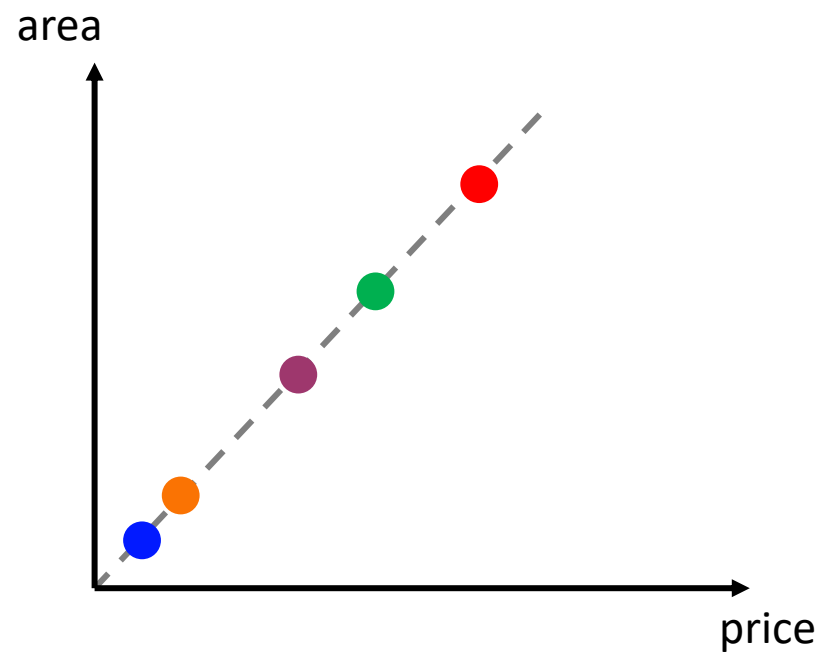
Standard algorithm compared to Strassen's algorithm, which uses one less scalar multiplication (7 instead of 8) for multiplying 2x2 matrices. Multiplications matter much more than additions for overall efficiency.



- a neural network architecture that incorporates problem-specific inductive biases
 - a procedure to generate useful synthetic data
 - a recipe to leverage symmetries of the problem.
- Traditional algorithm taught in school multiplies a 4x5 by 5x5 matrix using 100 multiplications;
 - this number was reduced to 80 with human ingenuity
 - AlphaTensor has found algorithms that do the same operation using just 76 multiplications.

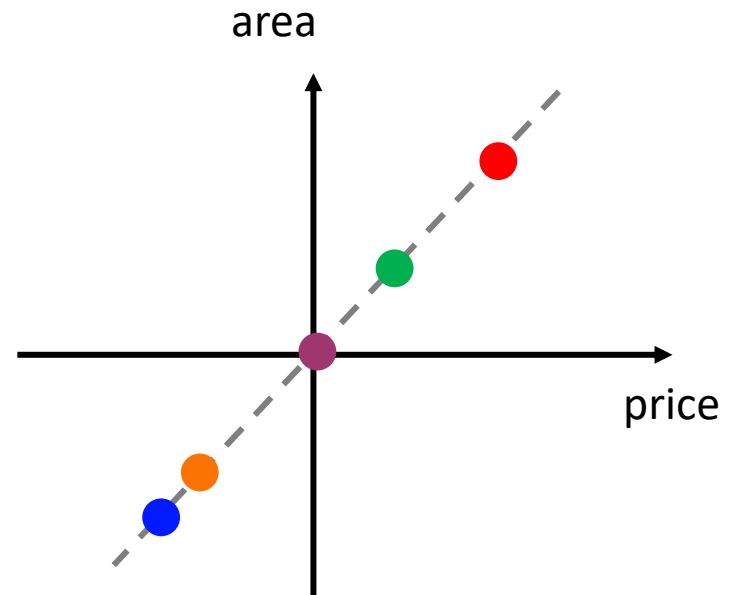
Idea of PCA

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

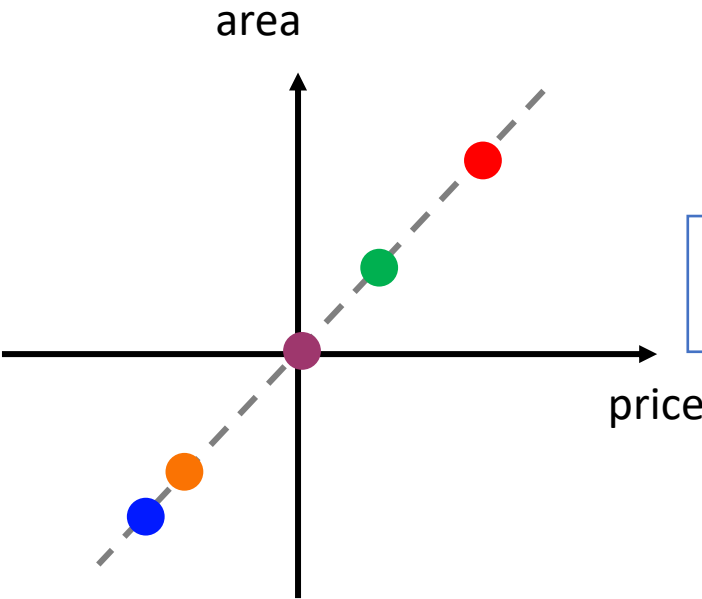


We subtract means from data points

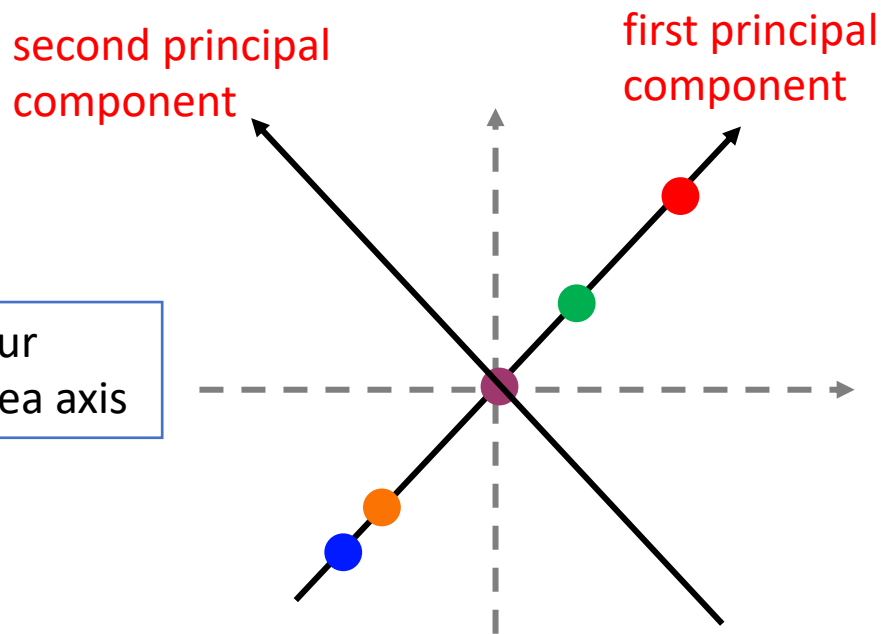
	House price (normalised)	House area (normalised)
a	5	5
b	-3	-3
c	2	2
d	-4	-4
e	0	0



Idea of PCA



We rotate our price and area axis

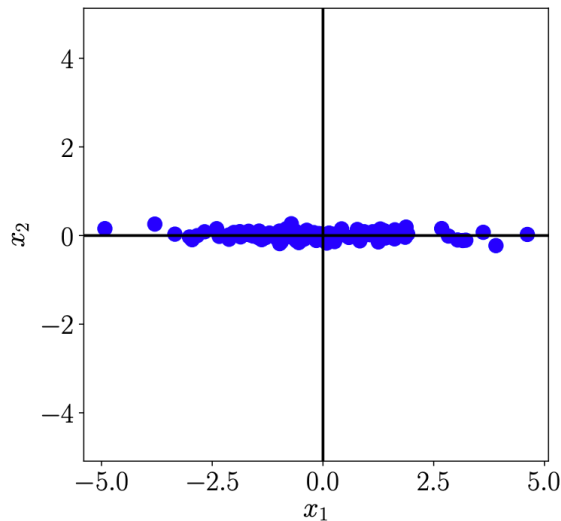


	House price (normalised)	House area (normalised)
a	5	5
b	-3	-3
c	2	2
d	-4	-4
e	0	0

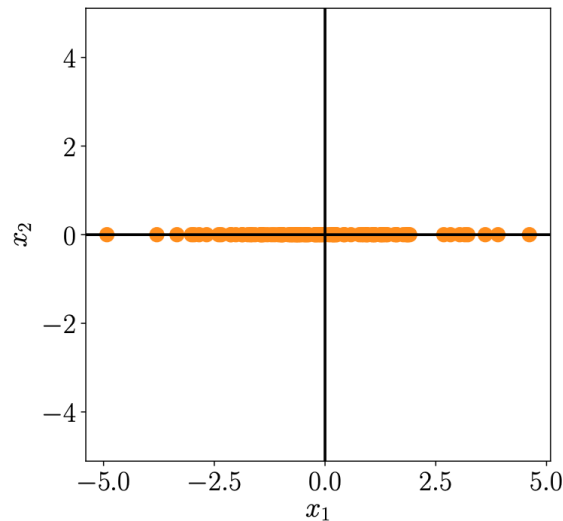
	First principal component	Second principal component
a	7.07	0
b	-4.24	0
c	2.82	0
d	-5.66	0
e	0	0

Motivation

- High-dimensional data, such as images, is hard to analyze, interpret, 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.



(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

- In PCA, we are interested in finding projections $\tilde{\mathbf{x}}_n$ 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
- 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**

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

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

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

of \mathbf{x}_n , where we define the projection matrix

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

- **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.,

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

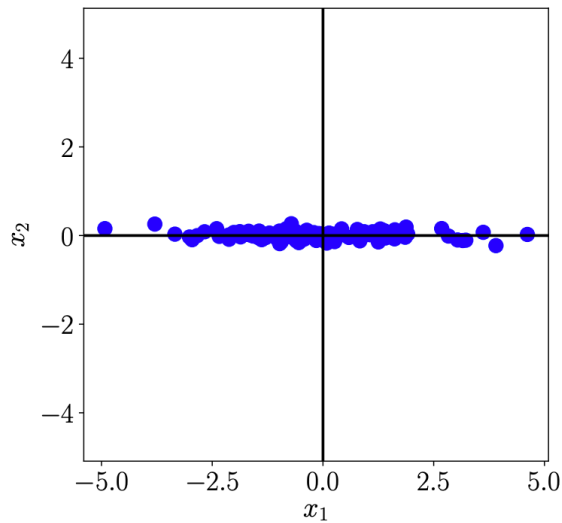
- However, when we consider vectors of the form

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

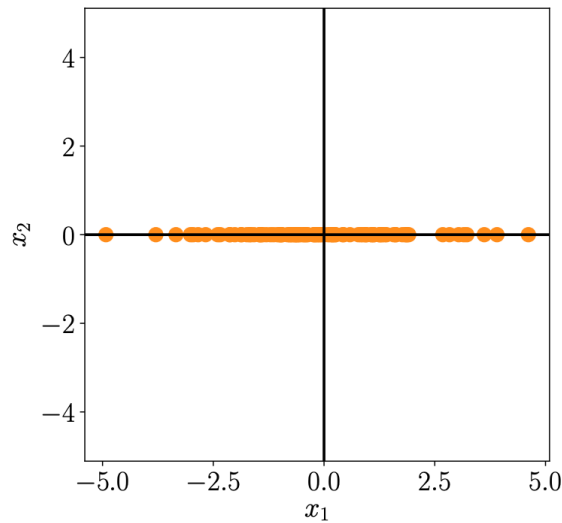
they can always be written as $0\mathbf{e}_1 + z\mathbf{e}_2$.

- 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



(a) Dataset with x_1 and x_2 coordinates.



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

- 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

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

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

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

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

$$\mathbb{V}[\mathbf{x} + \mathbf{y}] = \mathbb{V}[\mathbf{x}] + \mathbb{V}[\mathbf{y}] + \text{Cov}[\mathbf{x}, \mathbf{y}] + \text{Cov}[\mathbf{y}, \mathbf{x}]$$

- 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

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

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

- 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 λ_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 $\mathbf{0}$ gives us the relations

$$\begin{aligned} \mathbf{S} \mathbf{b}_1 &= \lambda_1 \mathbf{b}_1 \\ \mathbf{b}_1^T \mathbf{b}_1 &= 1 \end{aligned}$$

- We see that \mathbf{b}_1 is an eigenvector of \mathbf{S} , and λ_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

- Assume we have found the $m - 1$ eigenvectors of \mathbf{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 $\mathbf{b}_1 \cdots, \mathbf{b}_{m-1}$ from the data, and find principal components that compress the **remaining information**. We then arrive at the new data matrix,

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

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ 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 $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$.

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

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

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

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

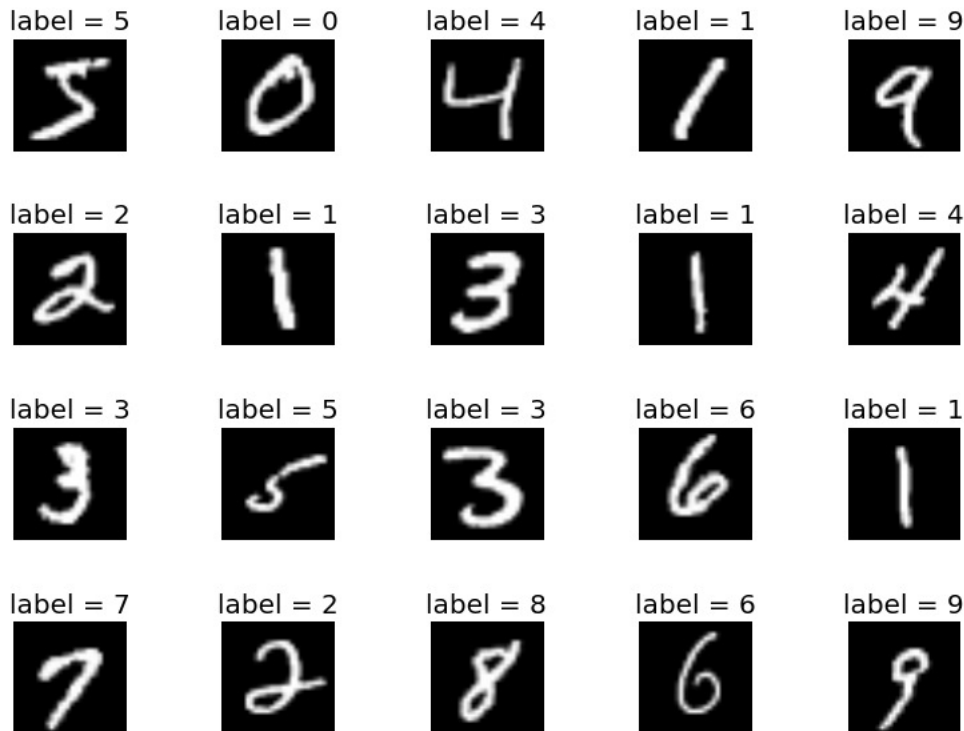
$$\hat{\mathbf{S}}\mathbf{b}_m = \mathbf{S}\mathbf{b}_m = \lambda_m\mathbf{b}_m \quad (1)$$

- \mathbf{b}_m is not only an eigenvector of \mathbf{S} but also of $\hat{\mathbf{S}}$.
- Specifically, λ_m is the largest eigenvalue of $\hat{\mathbf{S}}$ and λ_m is the m th largest eigenvalue of \mathbf{S} , and both have the associated eigenvector \mathbf{b}_m .
- Considering (1) and, $\mathbf{b}_m^T\mathbf{b}_m = 1$, the variance of the data projected onto the m th principal component is

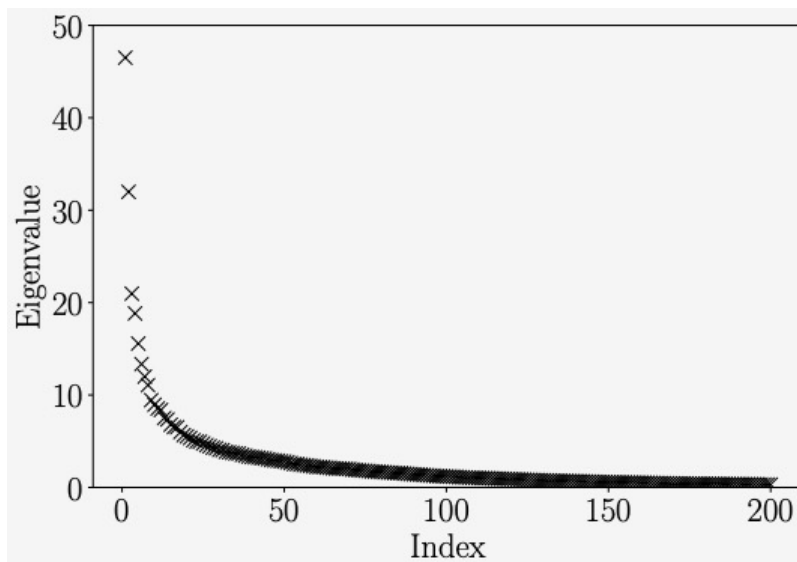
$$V_m = \mathbf{b}_m^T\mathbf{S}\mathbf{b}_m = \lambda_m\mathbf{b}_m^T\mathbf{b}_m = \lambda_m$$
- This means that the variance of the data, when projected onto an M -dimensional subspace, equals the sum of the eigenvalues that are associated with the corresponding eigenvectors of the data covariance matrix.

MNIST dataset

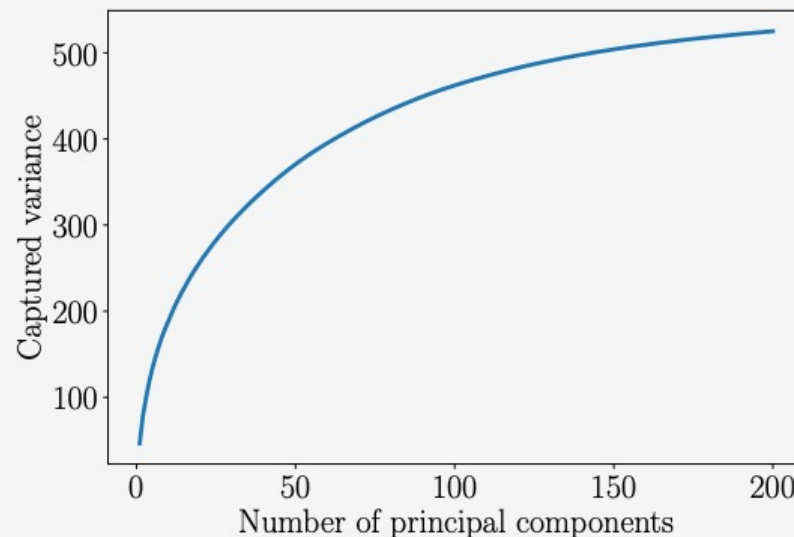
- 60,000 examples of handwritten digits 0 through 9.
- Each digit is a grayscale image of size 28×28 , i.e., it contains 784 pixels.
- We can interpret every image in this dataset as a vector $x \in \mathbb{R}^{784}$



Example - Eigenvalues of MNIST digit “8”



(a) Top 200 largest eigenvalues



(b) Variance captured by the principal components.

- 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

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

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

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

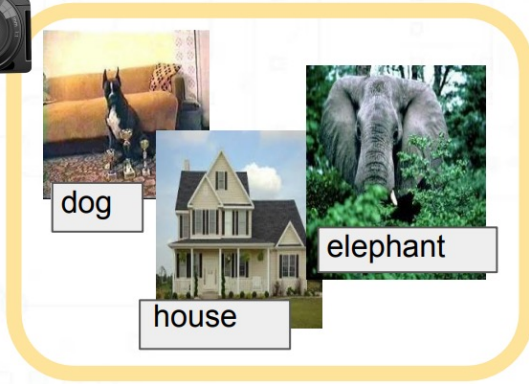
- The variance lost by data compression via PCA is

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

- 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 $1 - \frac{V_M}{V_D}$.

Domain Adaptation

Source Domain



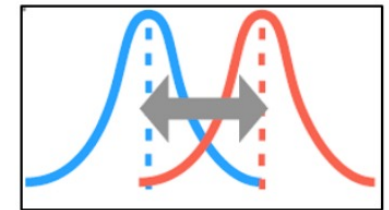
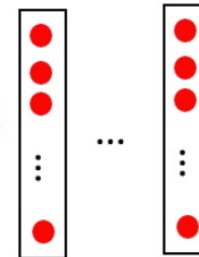
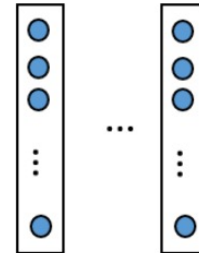
Target Domain



Source set

Source features

Source classifiers



Target set

Target features

Domain Alignment

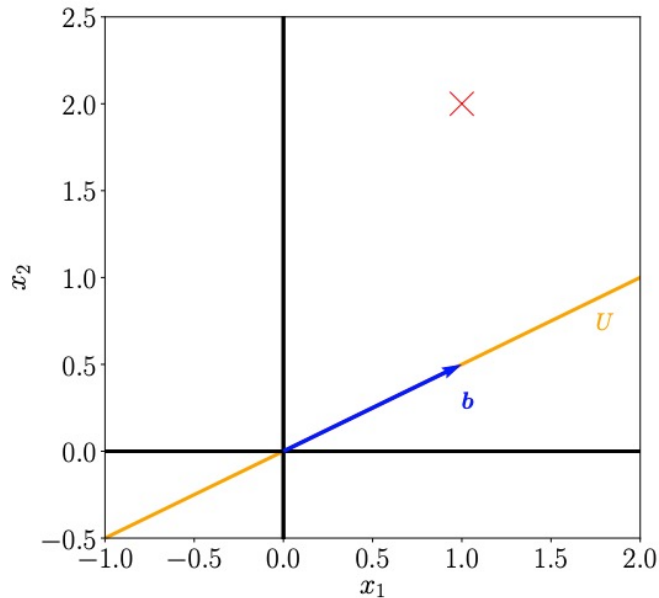
10.3 PCA from Projection Perspective

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

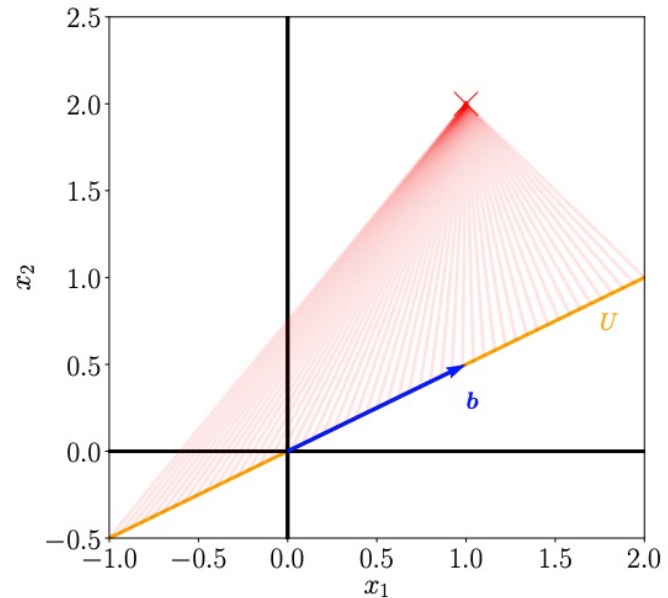
$$\begin{aligned} & \max_{\mathbf{b}_1} \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1 \\ & \text{subject to } \|\mathbf{b}_1\|^2 = 1 \end{aligned}$$

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

10.3.1 Setting and Objective



(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 $\|x - \tilde{x}\|$

- 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 :

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

for suitable coordinates $\zeta_d \in \mathbb{R}$.

- 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

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

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

- 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

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

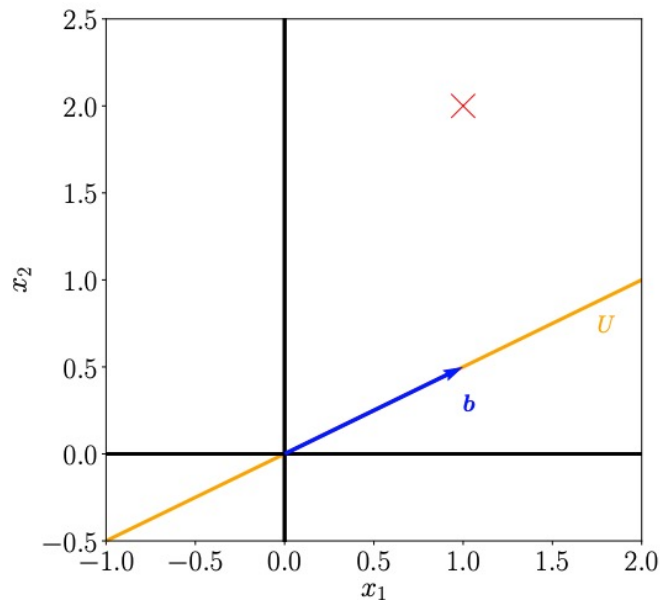
where $\mathbf{z}_n := [z_{1n}, \dots, z_{Mn}]^T \in \mathbb{R}^M$ is the coordinate vector of $\tilde{\mathbf{x}}_n$ with respect to the basis $(\mathbf{b}_1, \dots, \mathbf{b}_M)$.

- 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**)

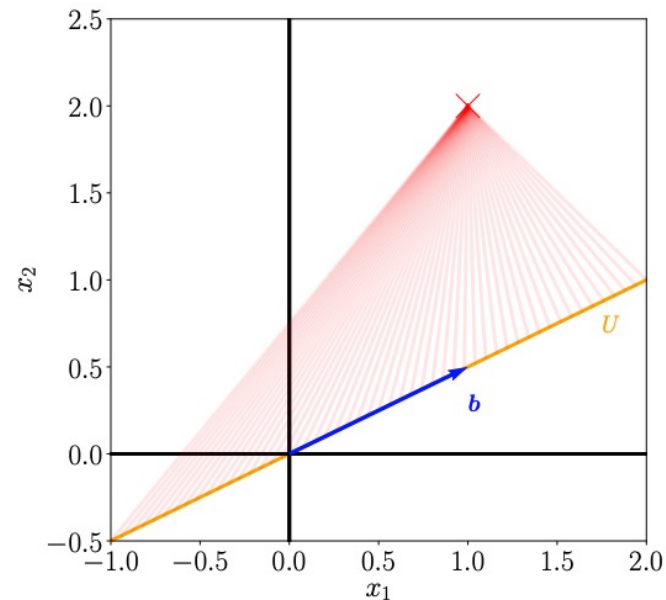
$$J_M := \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$$

- 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



(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

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

- 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

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

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

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

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

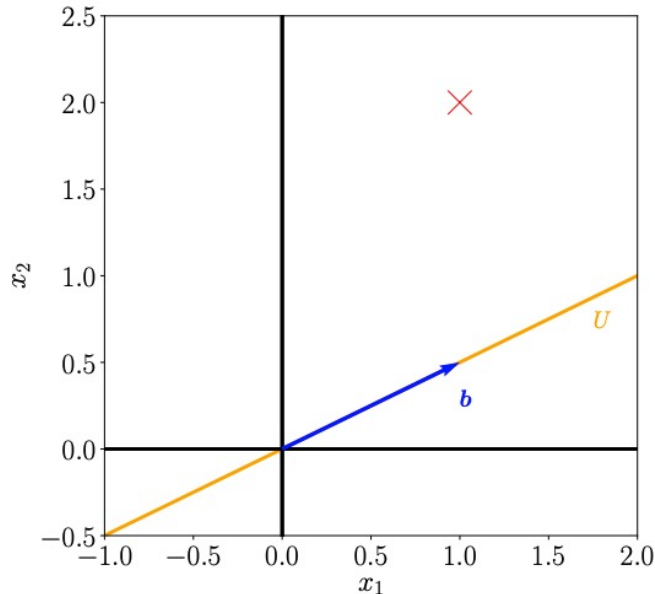
$$\frac{\partial J_M}{\partial z_{in}} = -\frac{2}{N}(\mathbf{x}_n^T \mathbf{b}_i - z_{in})$$

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

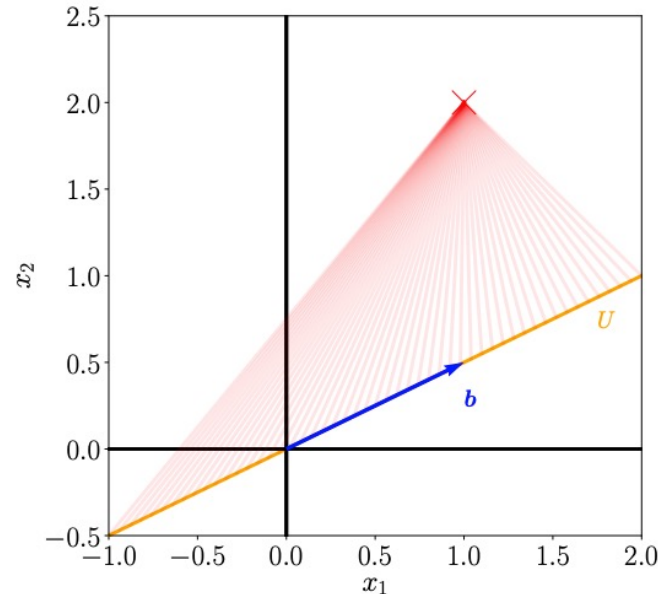
$$z_{in} = \mathbf{x}_n^T \mathbf{b}_i = \mathbf{b}_i^T \mathbf{x}_n$$

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

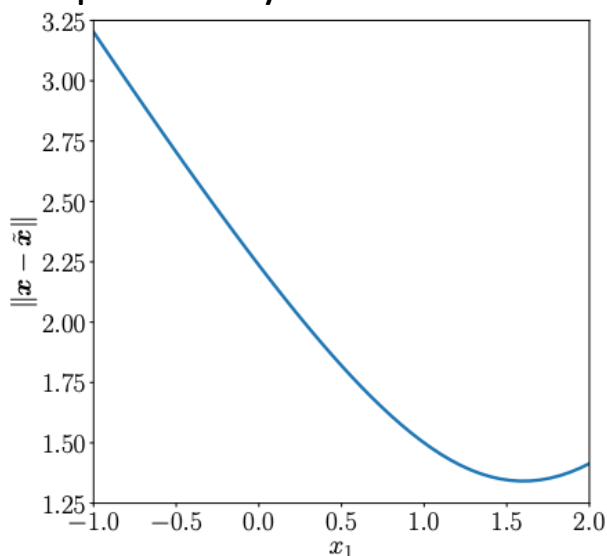
- The optimal coordinates z_{in} of the projection $\tilde{\mathbf{x}}_n$ are the coordinates of the orthogonal projection of the original data point \mathbf{x}_n onto the one-dimensional subspace that is spanned by \mathbf{b}_i .
- The optimal linear projection $\tilde{\mathbf{x}}_n$ of \mathbf{x}_n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



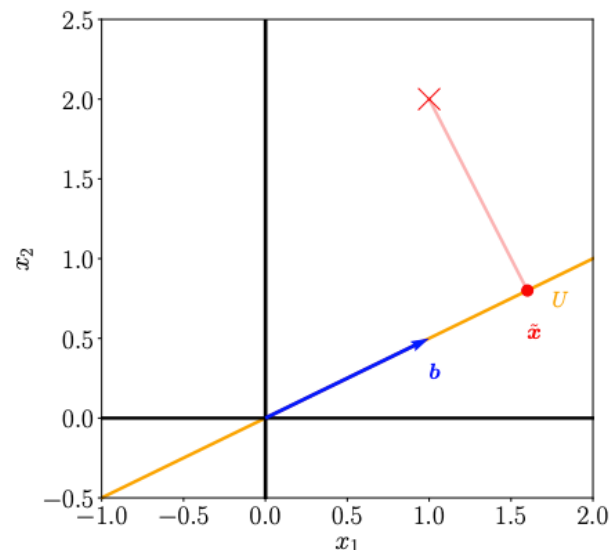
(a) A vector $x \in \mathbb{R}^2$ (red cross) shall be projected onto a one-dimensional subspace $U \subseteq \mathbb{R}^2$ spanned by b



(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 .

- 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 $z_j \mathbf{b}_j = \tilde{\mathbf{x}}$.

- 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

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

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

- 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

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

- We rearrange this equation, 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

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

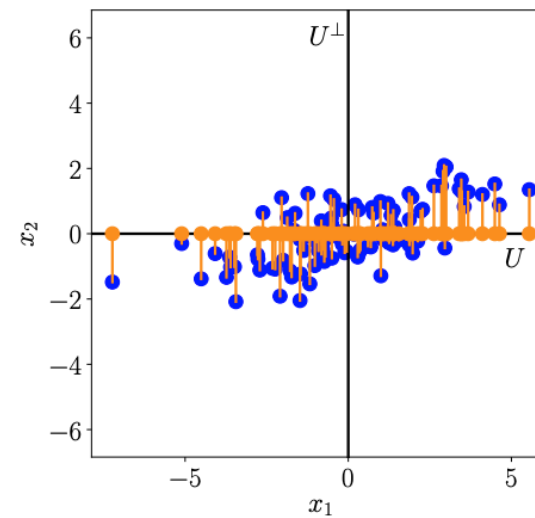
where we split the sum with D terms into a sum over M and a sum over $D - M$ terms.

- 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 .



- 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:

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

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

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

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

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

- 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 Slide #17, the average squared reconstruction error, when projecting onto the M -dimensional principal subspace, is

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

- where λ_j are the eigenvalues of the data covariance matrix.

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

- 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.

10.5 PCA in High Dimensions

- 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.

- 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.
- 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

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

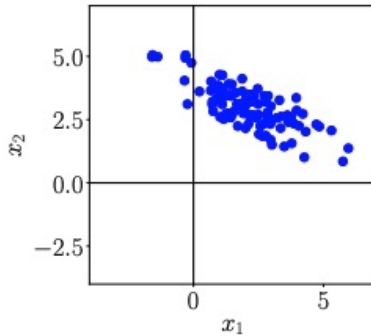
$$\frac{1}{N} \mathbf{X}^T \mathbf{X} \mathbf{c}_m = \lambda_m \mathbf{c}_m$$

- We get a new eigenvector/eigenvalue equation: λ_m remains eigenvalue.
- We obtain the eigenvector of the matrix $\frac{1}{N} \mathbf{X}^T \mathbf{X} \in \mathbb{R}^{N \times N}$ associated with λ_m as $\mathbf{c}_m := \mathbf{X}^T \mathbf{b}_m$.
- 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.
- 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

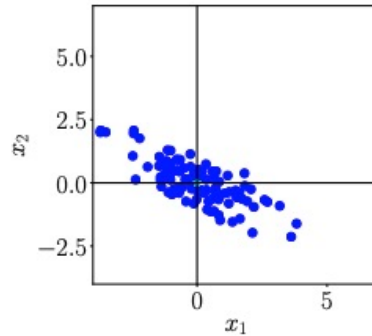
$$\underbrace{\frac{1}{N} \mathbf{X} \mathbf{X}^T}_{\mathbf{S}} \mathbf{X} \mathbf{c}_m = \lambda_m \mathbf{X} \mathbf{c}_m$$

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} .

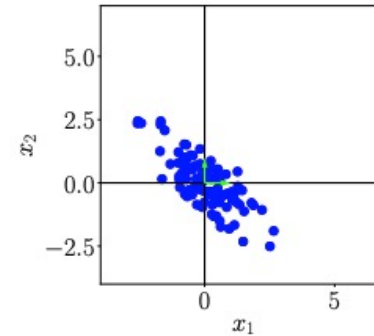
10.6 Key Steps of PCA in Practice



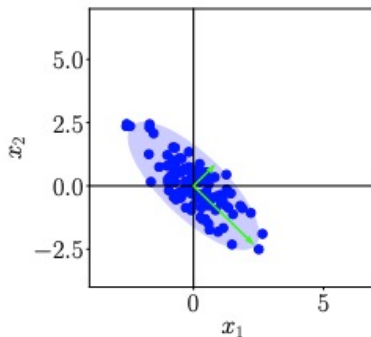
(a) Original dataset.



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

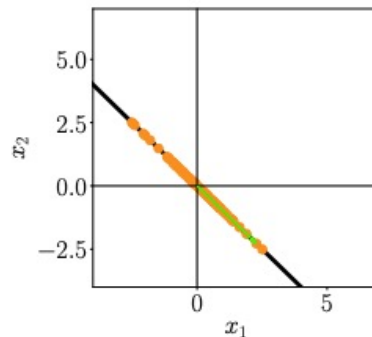


(c) Step 2: Dividing by the standard deviation to make the data unit free. Data has variance 1 along each axis.

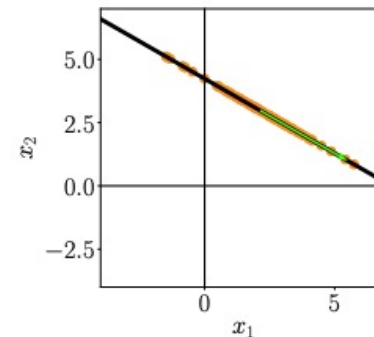


(d) Step 3: Compute eigenvalues and eigenvectors (arrows) of the data covariance matrix (ellipse).

eigendecomposition

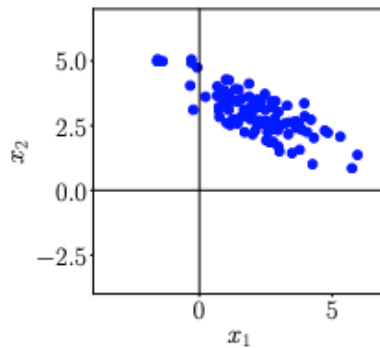


(e) Step 4: Project data onto the principal subspace.

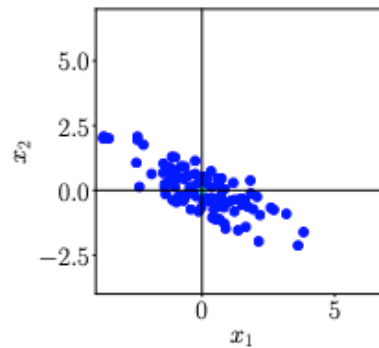


(f) Undo the standardization and move projected data back into the original data space from (a).

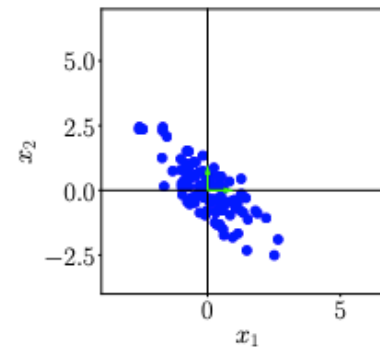
- **Step 1. Mean subtraction**
- We center the data by computing the mean μ of the dataset and subtracting it from every single data point. This ensures that the dataset has mean 0 .
- **Step 2. Standardization** Divide the data points by the standard deviation σ of the dataset for every dimension $d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



(a) Original dataset.

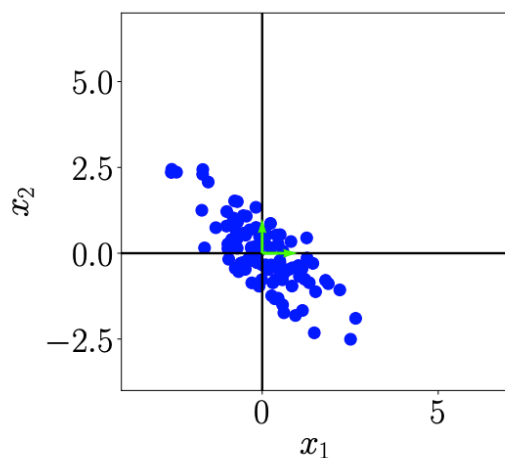


(b) Step 1: Centering by subtracting the mean from each data point.

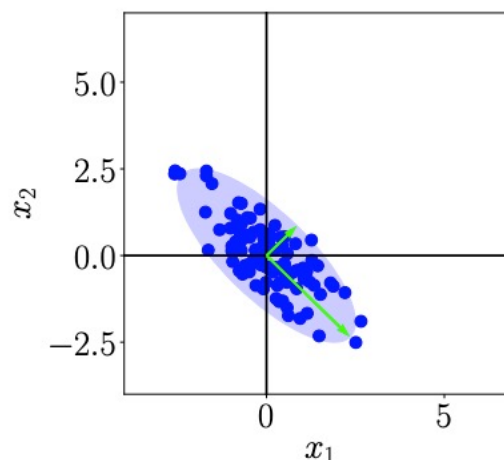


(c) Step 2: Dividing by the standard deviation to make the data unit free. Data has variance 1 along each axis.

- **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).

- **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 μ_d and standard deviation σ_d of the training data in the d th dimension, respectively, so that

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

where $x_*^{(d)}$ is the d th component of \mathbf{x}_* .

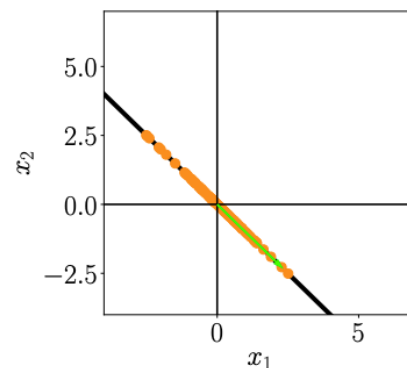
- We obtain the projection as

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

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.

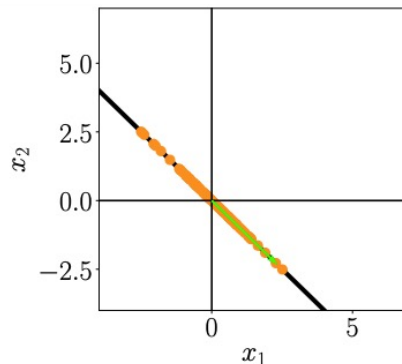


(e) Step 4: Project data onto the principal subspace.

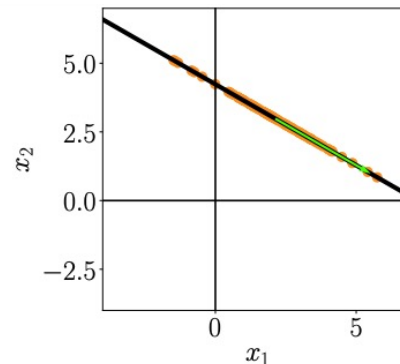
- Having standardized our dataset, $\tilde{\mathbf{x}}_* = \mathbf{B}\mathbf{B}^T\mathbf{x}_*$ only yields the projections in the context of the **standardized dataset**.
- To obtain our projection in the original data space (i.e., before standardization), we need to undo the standardization: multiply 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).