# Principal Component Analysis

Machine Learning 1 — Lecture 6
25[th] April 2023

Robert Peharz
Institute of Theoretical Computer Science
Graz University of Technology

- So far, we considered **linear models** for **supervised learning**
- That is, we predicted some **target** $y$ from an **input vector** $x$:
  - **linear regression** (for regression)
  - **logistic regression** (for classification)
  - both can be made non-linear, by replacing $x$ with a non-linear transformation $\phi(x)\colon \mathbb{R}^D \mapsto \mathbb{R}^K$
- Today, we will use a linear model for an **unsupervised learning** technique: **principal component analysis (PCA)**
- Thus, we will have a data set

$$\mathcal{D} = \left\{ x^{(1)}, x^{(2)}, \ldots, x^{(N)} \right\}$$

of input vectors, but **no target**
- Goal of unsupervised learning: find "interesting structure" in the data

PCA is a technique for **dimensionality reduction**.

Generally, dimensionality reduction means that we replace our data

$$\mathcal{D} = \left\{ \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots, \mathbf{x}^{(N)} \right\}$$

containing $D$-dimensional vectors with a new dataset

$$\mathcal{D}' = \left\{ \mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \ldots, \mathbf{z}^{(N)} \right\}$$

of $K$-dimensional vectors, $K \ll D$, where $\mathcal{D}'$ captures "certain characteristics" of $\mathcal{D}$.

Replace $\mathcal{D} = \left\{ \mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(N)} \right\}$ with $\mathcal{D}' = \left\{ \mathbf{z}^{(1)}, \ldots, \mathbf{z}^{(N)} \right\}$.
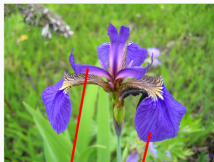
**Potential target characteristics (qualitative):**

- if $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$ are close, then also $\mathbf{z}^{(i)}$ and $\mathbf{z}^{(j)}$ should be close (**locality preserving**)
- $\mathbf{z}^{(i)}$ should allow for a good reconstruction of $\mathbf{x}^{(i)}$, i.e. there is some **decoder function** $g(\mathbf{z}^{(i)}) \approx \mathbf{x}^{(i)}$
- $\mathcal{D}'$ captures as much information or variance as possible of $\mathcal{D}$

**Applications of dimensionality reduction:**

- lossy compression
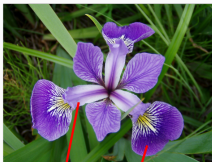- feature extraction (as input for other ML models)
- data visualization

**Iris Setosa**   **Iris Versicolor**   **Iris Virginica**

petal   sepal     petal   sepal     petal   sepal

Recall the **Iris dataset**, containing 4 features:

"sepallength", "sepalwidth", "petallength", "petalwidth", "class"
5.1, 3.5, 1.4, 0.2, Iris-setosa
4.9, 3.0, 1.4, 0.2, Iris-setosa
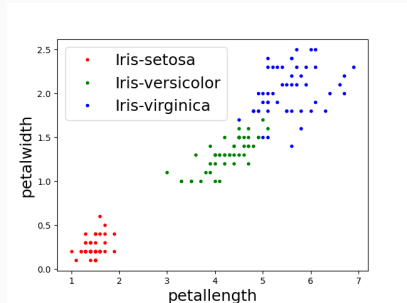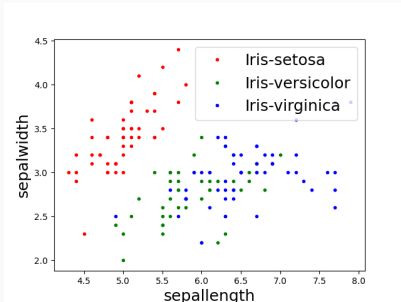. . .
7.0, 3.2, 4.7, 1.4, Iris-versicolor
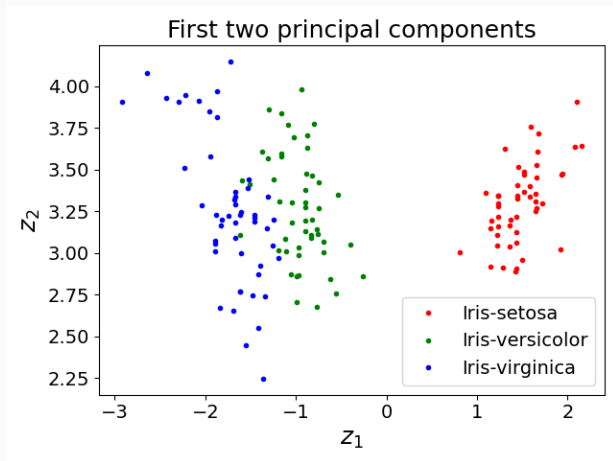6.4, 3.2, 4.5, 1.5, Iris-versicolor
. . .
6.3, 3.3, 6.0, 2.5, Iris-virginica
5.8, 2.7, 5.1, 1.9, Iris-virginica

## How to simultaneously plot 4 input features?

Computing the **first two principal components** transforms the data into a 2-dimensional space, which can be nicely visualized:



First two principal components
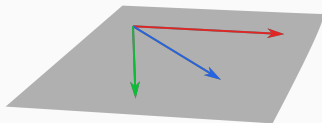
## Overview

Principal Component Analysis: Formulation

Computing the PCA

Example: MNIST Digits

Linear Discriminant Analysis
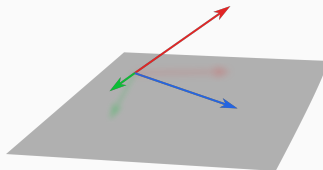
Vectors $x_1, x_2, \ldots, x_K$ are called **linearly independent**, when none of them can be expressed as a linear combination of the other vectors. That is, for each $1 \le i \le K$ and any coefficients $z_k$

$$x_i \neq \sum_{k=1, k \neq i}^{K} z_k x_k$$

**Linearly dependent vectors**    **Linearly independent vectors**

Let $B = (\boldsymbol{b}_1, \boldsymbol{b}_2, \ldots, \boldsymbol{b}_K)$ be a $D \times K$-matrix with orthonormal vectors as columns, where $K < D$. The **linear subspace** spanned by $B$ is defined as all possible linear combinations of $B$'s columns:

$$\boldsymbol{V} = \left\{ \boldsymbol{v} = B\boldsymbol{z} \mid \boldsymbol{z} = (z_1, \ldots, z_K)^T \in \mathbb{R}^K \right\}$$

Let $B = (\boldsymbol{b}_1, \boldsymbol{b}_2, \ldots, \boldsymbol{b}_K)$ be a matrix with $K$ orthonormal vectors as columns, spanning a $K$-dimensional subspace $\boldsymbol{V}$ of $\mathbb{R}^D$.

We can **project** an arbitrary vector $\boldsymbol{x} \in \mathbb{R}^D$ onto $\boldsymbol{V}$ by

- computing the **projection coefficients** $B^T \boldsymbol{x} =: \boldsymbol{z} \in \mathbb{R}^K$

- computing the **projection**/**reconstruction** $B\boldsymbol{z} =: \hat{\boldsymbol{x}} \in \boldsymbol{V}$

- thus, $\hat{\boldsymbol{x}} = \underbrace{BB^T}_{\text{projection matrix}} \boldsymbol{x}$

- $\hat{\boldsymbol{x}}$ is the **closest point** in $\boldsymbol{V}$ (in Euclidean distance) to $\boldsymbol{x}$

- the **residual** $\boldsymbol{r} = \boldsymbol{x} - \hat{\boldsymbol{x}}$ is always orthogonal to $\hat{\boldsymbol{x}}$

Idea of PCA: Learn an "interesting" subspace $V$.

Equivalently, learn a $D \times K$-matrix $B$ of orthonormal vectors capturing $K$ "interesting directions."
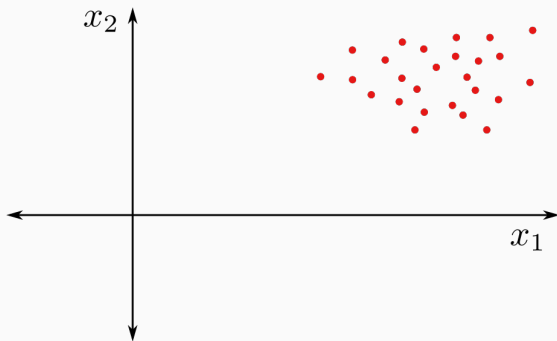
# Principal Component Analysis: Formulation

## Data

PCA is an **unsupervised** technique, i.e. we have a training dataset of feature vectors

$$\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots, \mathbf{x}^{(N)}\}$$

We can interpret the data as a "point cloud" in $\mathbb{R}^D$:

## Step 1: Centering the Data

For PCA we assume that the average of the data is $\bar{x} = \mathbf{0}$:

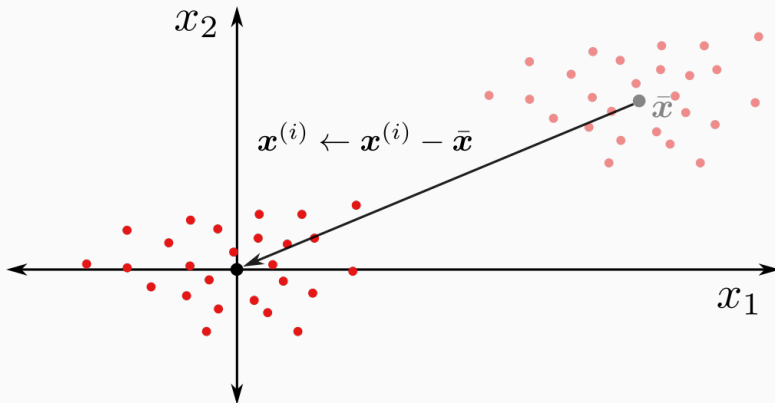$$\bar{x} := \frac{1}{N} \sum_{i=1}^{N} x^{(i)} = \mathbf{0}$$

If not, we simply remove the average from each training example:

$$x^{(i)} \leftarrow x^{(i)} - \bar{x}$$

This effectively moves the center of the point cloud to the origin.

Thus, we will w.l.o.g. assume that $\bar{\mathbf{x}} = \mathbf{0}$.

## Goal of PCA

- Find an "interesting subspace" $V$ of the data
- This amounts to finding a matrix $B = (\boldsymbol{b}_1, \ldots, \boldsymbol{b}_K)$ of $K$ orthonormal basis vectors
- The basis vectors $\boldsymbol{b}_k$ define $K$ "interesting directions"
- $K$ is a **hyper-parameter** selected by the user
- For now, let's consider the special case $K = 1$, i.e. we want to find the "most interesting direction" in the data
- How to define "interestingness"?

Which of the three directions $b$, $b'$, $b''$ is the most "interesting" one? Why?

$\boldsymbol{b'}$ is the direction where

1. the **variance of the projected data is maximal**
2. the **sum of squared projection errors is minimal**

Let $\mathcal{D} = \{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(N)}\}$ be a centered training set of $D$-dimensional vectors.

**Among all unit vectors,** let $\boldsymbol{b}_1$ be the unit vector where

- the **variance of the projected data**
  $\textbf{var}\left(\left\{ z_1^{(1)}, z_1^{(2)}, \ldots, z_1^{(N)} \right\}\right)$ is **maximal**, where

$$z_1^{(i)} = \boldsymbol{b}_1^T \boldsymbol{x}^{(i)}$$

  or, **equivalently,** where

- the sum of **squared projection errors**

$$\sum_{i=1}^{N} \|\boldsymbol{x}^{(i)} - \hat{\boldsymbol{x}}^{(i)}\|_2^2 = \sum_{i=1}^{N} \|\boldsymbol{x}^{(i)} - \boldsymbol{b}_1 z_1^{(i)}\|_2^2$$

  is **minimal.**

The vector $\boldsymbol{b}_1$ is called the **first principal direction** of dataset $\mathcal{D}$.

The values $\left\{ z_1^{(1)}, z_1^{(2)}, \ldots, z_1^{(N)} \right\}$ are called the **first principal component** (**scores**) of $\mathcal{D}$.

- the scores $z_1^{(1)}, z_1^{(2)}, \ldots, z_1^{(N)}$ are a 1-dimensional representation of the original samples $\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots, \boldsymbol{x}^{(N)}$

- the vector $\hat{\boldsymbol{x}}^{(i)} = \boldsymbol{b}_1 z_1^{(i)} = \boldsymbol{b}_1 \boldsymbol{b}_1^T \boldsymbol{x}^{(i)}$ is the projection of $\boldsymbol{x}^{(i)}$ onto $\boldsymbol{b}_1$ (or, the best reconstruction of $\boldsymbol{x}^{(i)}$ given $z_1^{(i)}$)

- $\boldsymbol{b}_1$ is the **direction in which the data varies the most**, or the **direction of highest energy** – often, although not always, this direction carries the "most useful information" in the data

- the equivalence between **maximal variance** and **minimal projection/reconstruction error** (not proven here) is an example of **duality**, which occurs often in optimization

- by subtracting the projections $\hat{x}^{(i)}$ from the original samples $x^{(i)}$ and finding the principal component in the **residual vectors** $r^{(i)} = x^{(i)} - \hat{x}^{(i)}$, one finds the **second principal component**, and so forth

- we will see in the following, however, that is is actually quite simple to find the first $K$ principal components in one step

# Computing the PCA

- Assume that the dataset $\mathcal{D} = \{x^{(1)}, \ldots, x^{(N)}\}$ is already centered, i.e. $\bar{x} = 0$
- We want to find the **first principal direction** $b_1$
- To simplify notation, we denote it by $b$, dropping the subscript

## Problem Setup

- The PCA scores are given as

$$z^{(i)} = \boldsymbol{b}^T \boldsymbol{x}^{(i)}$$

- We wish to maximize the **empirical variance**

$$\textbf{var}(z^{(i)}) = \frac{1}{N} \sum_{i=1}^{N} (z^{(i)})^2 = \frac{1}{N} \sum_{i=1}^{N} (\boldsymbol{b}^T \boldsymbol{x}^{(i)})(\boldsymbol{x}^{(i)}{}^T \boldsymbol{b})$$

- We require that $\boldsymbol{b}$ is a **unit vector**, i.e.

$$\boldsymbol{b}^T \boldsymbol{b} = 1, \qquad \sum_{i=1}^{D} b_i^2 = 1$$

(note that if $\boldsymbol{b}$ was not constrained to be a unit vector, $\textbf{var}(z^{(i)})$) would grow to infinity with the norm of $\boldsymbol{b}$)

We can transform the problem as follows:

$$
\begin{aligned}
\max_{\boldsymbol{b}^T \boldsymbol{b}=1} \ \mathbf{var}(z^{(i)}) &= \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{b}^T \boldsymbol{x}^{(i)} \boldsymbol{x}^{(i)^T} \boldsymbol{b} \\
&= \boldsymbol{b}^T \underbrace{\left( \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}^{(i)} \boldsymbol{x}^{(i)^T} \right)}_{=:\textbf{covariance matrix } C} \boldsymbol{b}
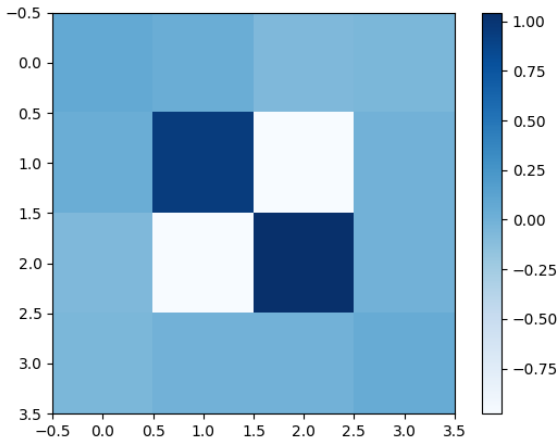\end{aligned}
$$

## Covariance Matrix

The (empirical) **covariance matrix** $C$ is a $D \times D$ matrix containing all possible empirical (co-)variances in the data:

$$C = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}^{(i)} \boldsymbol{x}^{(i)^T} = \frac{1}{N} \boldsymbol{X}^T \boldsymbol{X}.$$

(**Recall** that $\boldsymbol{X}$ is the design matrix, whose rows are the $N$ training examples).

- the diagonal element $C_{dd}$ is the **empirical variance** of the $d^{\text{th}}$ dimension: $\frac{1}{N} \sum_i x_d^{(i)} x_d^{(i)}$
- the off-diagonal element $C_{de}$ is the **empirical co-variance** between $d^{\text{th}}$ and $e^{\text{th}}$ dimension: $\frac{1}{N} \sum_i x_d^{(i)} x_e^{(i)}$
- $C$ is always a symmetric matrix

## PCA: Optimization Problem

With the covariance matrix, the PCA problem can compactly be written as

$$\max_{\boldsymbol{b}^T \boldsymbol{b}=1} \boldsymbol{b}^T C \boldsymbol{b}$$

What is the optimal $\boldsymbol{b}$ here? Evidently, $\boldsymbol{b}$ must be a vector somehow connected with $C$.

- Key and solution to the PCA problem: **eigen decomposition** of $C$

- For any matrix $C$, when it holds for some vector $\mathbf{v}$ and a scalar $\lambda$ that

$$C\mathbf{v} = \lambda\mathbf{v}$$

  then we call $\mathbf{v}$ an **eigen vector** of $C$ and $\lambda$ its corresponding **eigen value**

- In this case, also $\alpha\mathbf{v}$ is an eigen vector for any $\alpha \neq 0$; thus, we consider only normalized eigen vectors

For any **symmetric** $D \times D$-matrix $C$, there are $D$ **orthonormal** eigen vectors $\boldsymbol{v}_1, \boldsymbol{v}_2, \ldots, \boldsymbol{v}_D$ and $D$ sorted eigen values $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_D$, such that $C$ can be written as

$$C = VEV^T$$

where $V = (\boldsymbol{v}_1, \boldsymbol{v}_2, \ldots, \boldsymbol{v}_D)$ is an **orthonormal matrix** and $E$ is diagonal:

$$E = \begin{pmatrix} \lambda_1 & 0 & \ldots & 0 \\ 0 & \lambda_2 & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \ldots & \lambda_D \end{pmatrix}$$

For covariance matrices it further holds that all $\lambda_i \geq 0$; they are so-called **positive definite matrices**.

## PCA Solution

Using the eigen decomposition we can write the PCA problem as

$$\max_{\boldsymbol{b}^T \boldsymbol{b}=1} \boldsymbol{b}^T C \boldsymbol{b} = \underbrace{\boldsymbol{b}^T V}_{\boldsymbol{a}^T} \, E \, \underbrace{V^T \boldsymbol{b}}_{\boldsymbol{a}} = \sum_{i=1}^{D} \lambda_i a_i^2$$

- Since $\boldsymbol{b}$ is a unit vector and $V$ is orthonormal, also $\boldsymbol{a}$ is a unit vector, i.e. $\sum_{i=1}^{D} a_i^2 = 1$
- Evidently, $a_i^2 = (\boldsymbol{v}_i^T \boldsymbol{b})^2 \geq 0$
- Thus, the values $a_i^2$ are **non-negative** and **sum to one**
- The variance is

$$\lambda_1 a_1^2 + \lambda_2 a_2^2 + \cdots + \lambda_D a_D^2$$

- When is this maximal?

## PCA Solution cont'd

- The variance becomes maximal when $a_1^2 = 1$:

$$\lambda_1 1 + \lambda_2 0 + \cdots + \lambda_D 0$$

- Since $a_i = (\boldsymbol{v}_i^T \boldsymbol{b})$, this exactly happens if $\boldsymbol{b} = \boldsymbol{v}_1$!

- **Thus, the first eigen vector $\boldsymbol{v}_1$ of the covariance matrix is the first principal direction $\boldsymbol{b}_1$**

- The corresponding eigen value $\lambda_1$ is the variance which is "captured" in this direction (**explained variance**)
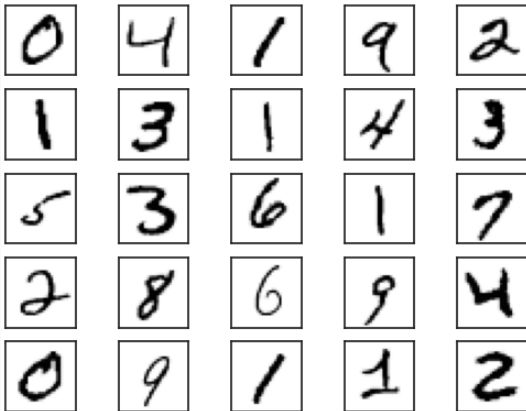
## PCA Solution cont'd

- More generally, the first $K$ eigen vectors $\mathbf{v}_1, \mathbf{v}_2, \ldots, \mathbf{v}_K$ are the the **first $K$ principal directions**

- The sum of their corresponding eigen values is the variance explained by the first $K$ principal components:

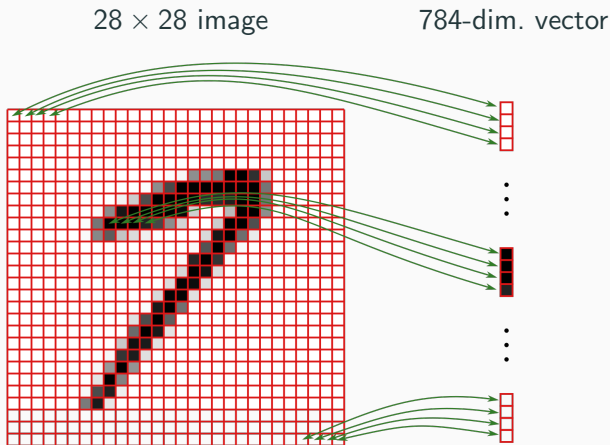$$\text{"explained variance"} = \sum_{i=1}^{K} \lambda_i$$

# Example: MNIST Digits

- collection of 60,000 images of handwritten digits
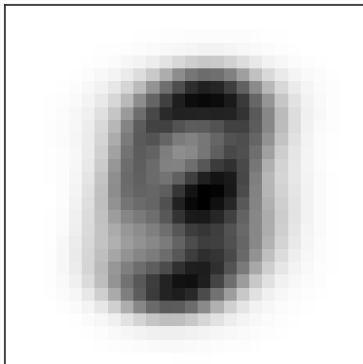- images are $28 \times 28$ pixels, with 256 gray scales
- white $= 0$, black $= 255$

By following a fixed pixel order (e.g. row-wise) we can convert each image into a vector. By "uprolling" the vector, we get back an image.
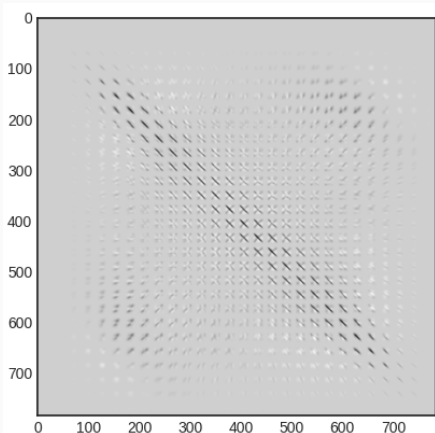
$28 \times 28$ image          784-dim. vector

The first step is the center the data, i.e. subtract the average vector (image) from each sample:

$$\bar{x} = \frac{1}{60,000} \sum_{i=1}^{60,000} x^{(i)}, \qquad \forall i: \ x^{(i)} \leftarrow x^{(i)} - \bar{x}$$

$\bar{x} =$

$$C = \frac{1}{60,000} \sum_{i=1}^{60,000} \mathbf{x}^{(i)} \mathbf{x}^{(i)^T} = \frac{1}{N} \mathbf{X}^T \mathbf{X}$$

$$V, E = \textbf{eigen}(C)$$

The first $K$ columns of $V$ contain the first $K$ **principal directions**.
We store them in a $784 \times K$ matrix $B$.

**Note:** Some implementations return the eigenvalues in arbitrary order.
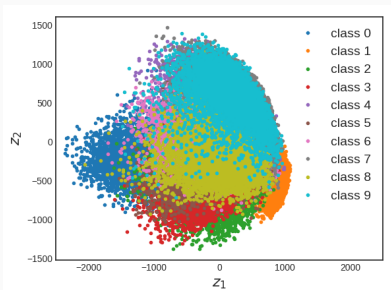In this case, sort them in descending order and sort the columns of $V$
correspondingly.

black pixels correspond to large positive values, white pixels correspond to small negative values
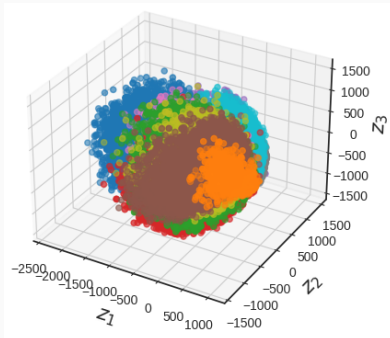
The matrix $B$ projects the data into a $K$-dimensional subspace:

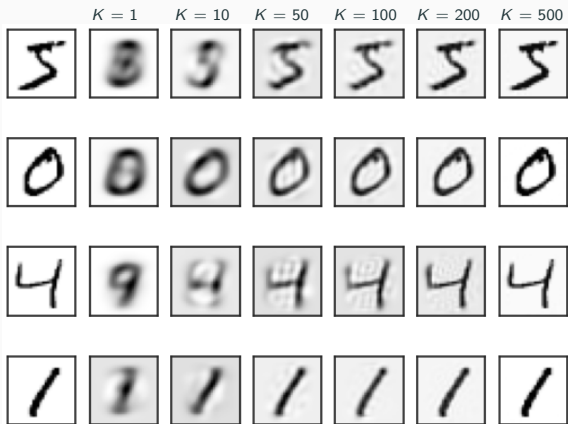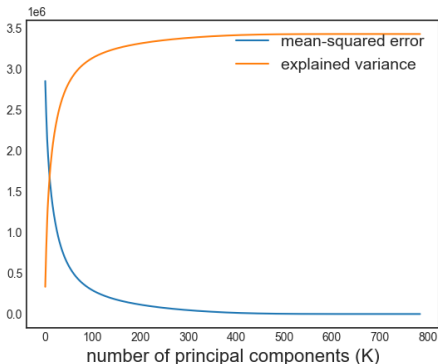$$\boldsymbol{Z} = \boldsymbol{X}B$$

$K = 2$                    $K = 3$

Reconstructions of four samples (first column) with 1, 10, 50, 100, 200 and 500 PCA components:

Reconstruction error (measured in mean squared error) and explained variance ($\sum_i \lambda_i$) over $K$:



$\rightarrow$ can be used to select $K$

**Given:** dataset $\mathcal{D} = \{\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots, \boldsymbol{x}^{(N)}\}$, parameter $K$

1. Center data
   - $\bar{\boldsymbol{x}} = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}^{(i)}$
   - for all $i$: $\boldsymbol{x}^{(i)} \leftarrow \boldsymbol{x}^{(i)} - \bar{\boldsymbol{x}}$
2. Compute covariance matrix
   $C = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}^{(i)} \boldsymbol{x}^{(i)^T} = \frac{1}{N} \boldsymbol{X}^T \boldsymbol{X}$
3. Compute eigen decomposition: $V, E \leftarrow$ **eigen**$(C)$
4. Let $B$ be the matrix consisting of the first $K$ columns of $V$
5. **return** $B$

   - The **principal component scores** (**low-dimensional features**) are then given as $\boldsymbol{Z} = \boldsymbol{X}B$ ($N \times K$-matrix)
   - Projections/reconstructions are given as $\hat{\boldsymbol{X}} = \boldsymbol{Z}B^T$
   - Add $\bar{\boldsymbol{x}}$ again to the reconstructions
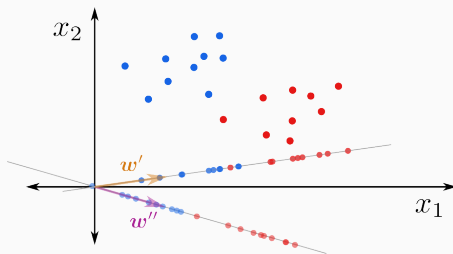
# Linear Discriminant Analysis

## Supervised Dimensionality Reduction

- PCA is **unsupervised**, i.e. there are no targets
- If one additionally has target values, one might use them to guide dimensionality reduction
- A classical method is Fisher's **linear discriminant analysis (LDA)**, using class information

## Separating Classes

- Let $\mathcal{D} = \left\{(\mathbf{x}^{(i)}, y^{(i)})\right\}_{i=1}^{N}$ be a binary classification dataset, i.e. each $y^{(i)} \in \{-1, 1\}$
- We want to learn a linear function $z = f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ projecting $D$-dimensional vectors onto the real line
- This function should allow us to separate the two classes well
- How to select $\mathbf{w}$?



**Note:** $\|\mathbf{w}\|$ is an arbitrary scaling, hence we can assume $\|\mathbf{w}\| = 1$.

## Separating the Means?

**First idea:** Maximize the distance between projected means

- Let $\boldsymbol{m}^-$ and $\boldsymbol{m}^+$ be the means of the negative and positive classes, respectively:

$$\boldsymbol{m}^- = \frac{1}{N^-} \sum_{i:\, y^{(i)}=-1} \boldsymbol{x}^{(i)} \qquad \boldsymbol{m}^+ = \frac{1}{N^+} \sum_{i:\, y^{(i)}=+1} \boldsymbol{x}^{(i)}$$
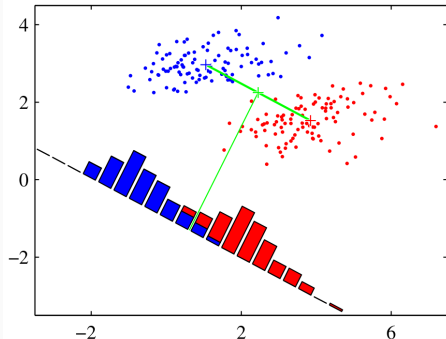
  Here $N^-$ and $N^+$ are the number of positive and negative samples, respectively.

- Let $m^-$ and $m^+$ be the projected means, i.e.

$$m^- = \boldsymbol{w}^T \boldsymbol{m}^- \qquad m^+ = \boldsymbol{w}^T \boldsymbol{m}^+$$

- It can be shown, that $\boldsymbol{w} \propto \boldsymbol{m}^- - \boldsymbol{m}^+$, i.e. the vector pointing from one mean to the other, maximizes the distance between the projected means $m^-$ and $m^+$

## Separating the Means?



- The idea works to a certain extent, but apparently we are not doing a perfect job
- Considerable overlap between the classes in projected space
- The problem is that we ignored the **covariance structure** of the point clouds

- Let $z^{(i)} = \boldsymbol{w}^T \boldsymbol{x}^{(i)}$ be the $i^{\text{th}}$ projected sample
- Let again $m^-$ and $m^+$ be the projected means
- Further, let $v^-$ and $v^+$ be the **within-class variances** of the projected data:

$$v^- = \frac{1}{N^-} \sum_{i\,:\,y^{(i)}=-1} \left(z^{(i)} - m^-\right)^2 \quad v^+ = \frac{1}{N^+} \sum_{i\,:\,y^{(i)}=+1} \left(z^{(i)} - m^+\right)^2$$

- The **Fisher criterion** is

$$J(\boldsymbol{w}) = \frac{(m^- - m^+)^2}{v^- + v^+}$$

- Increasing $J(\boldsymbol{w})$
  - increases the quadratic distance between $m^-$ and $m^+$
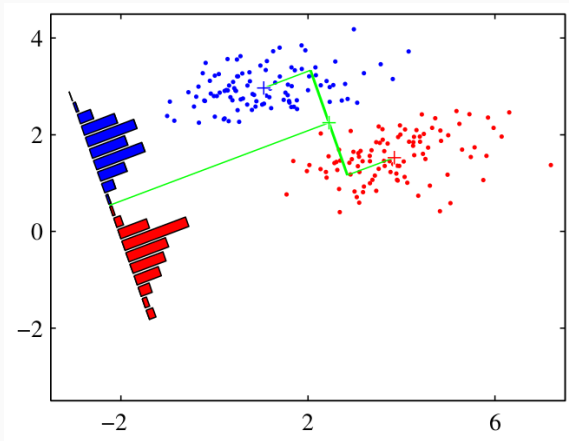  - and/or decreases variance (spread) within each class

- Fisher criterion $J(\mathbf{w})$ has an analytic maximum, denoted as **linear discriminant analysis (LDA)**

- Let $C^-$ and $C^+$ be the covariance matrices of the negative and positive class in the original space

$$C^- = \frac{1}{N^-} \sum_{i\,:\,y^{(i)}=-1} (\mathbf{x}^{(i)} - \mathbf{m}^-)(\mathbf{x}^{(i)} - \mathbf{m}^-)^T$$

$$C^+ = \frac{1}{N^+} \sum_{i\,:\,y^{(i)}=+1} (\mathbf{x}^{(i)} - \mathbf{m}^+)(\mathbf{x}^{(i)} - \mathbf{m}^+)^T$$

- Then the maximum of $J(\mathbf{w})$ is

$$\mathbf{w}_{LDA} \propto (C^- + C^+)^{-1}(\mathbf{m}^- - \mathbf{m}^+)$$

In contrast to the previous solution, LDA perfectly separates the two classes in this example:

- Similar to PCA, LDA uses covariance structure to determine a linear subspace

- Main difference is that PCA is unsupervised, while LDA is supervised

- LDA can be generalized to many classes as well, yielding a $(|\mathcal{C}| - 1)$-dimensional subspace