# Removing Redundancies From Data: Principle Component Analysis

COMS21202, Part III

- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

1

---

## Objectives

- Understand potential harm of high dimensionality of dataset
- Use Principle Component Analysis (PCA) to remove "redundant" dimensions from data.

- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

2

---

## High Dimensionality, Good? Bad?

- $X = \{x_i\}_{i=1}^n, x \in R^d$.
- Is a large $d$ always a good thing?
  - ☺ We have more info as $d$ grows!
  - ☹ LS does not work when $d > n$
  - ☹ Large $d$ causes overfitting
  - More ☹ ?

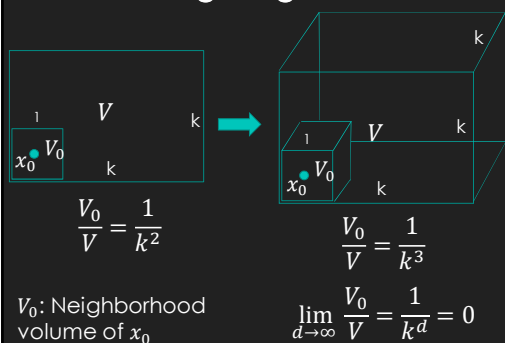- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

3

## Curse of Dimensionality (CoD)

- CoD is a generic term referring to the fact that many machine learning algorithms scale very poorly with $d$, in terms of performance.
  - Many geometry concepts work differently in higher dimensional space.
  - One of those concepts is "locality".

- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

4

## The Vanishing Neighborhood

$$\frac{V_0}{V} = \frac{1}{k^2}$$

$$\frac{V_0}{V} = \frac{1}{k^3}$$

$V_0$: Neighborhood volume of $x_0$

$$\lim_{d\to\infty} \frac{V_0}{V} = \frac{1}{k^d} = 0$$

- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

5

## The Vanishing Neighborhood

- The neighborhood cube quickly vanishes as $d$ increases.
- As a result, your k-nearest neighbors are **no longer** in the neighborhood $V_0$.
- These neighbors are no longer good at predicting the label of $x_0$.

- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

6

○ Can we reduce the dimensionality of $X$ without losing too much information?

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

7

## Reduce the Dimensionality using Feature Transform

○ We want to find a **feature transform** $f(x) \in R^m$, where $m \ll d$.

  ○ $f$ transforms original input $x$ to a subspace as $R^m \subset R^d$.

  ○ We assume our dataset is **centered**:

    ○ $\frac{1}{n}\sum_{i=1}^{n} x_i = 0$

    ○ **If dataset $X'$ is not centered:**

    ○ **Centering:** $\forall_i \, x_i = x'_i - \frac{1}{n}\sum_{i=1}^{n} x'_i$

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.
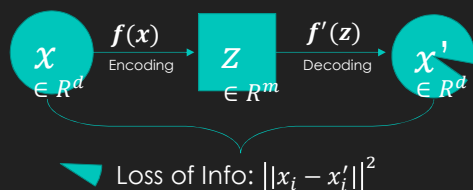
8

## Reduce the Dimensionality using Feature Transform

○ What is the optimal strategy of selecting $f(x)$?

○ Want to reduce dimension using $f$.

  ○ while preserving **as much info as possible**!

○ Let's look at this problem from data compression perspective!

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.
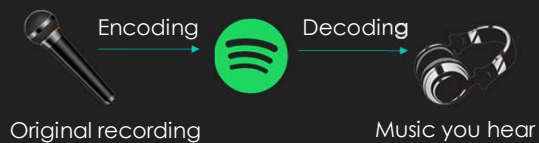
9

## Encoder and Decoder



$$x \in R^d \xrightarrow[\text{Encoding}]{f(x)} z \in R^m \xrightarrow[\text{Decoding}]{f'(z)} x' \in R^d$$

Loss of Info: $\left\|x_i - x_i'\right\|^2$

- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

10

## Codec



Encoding → Decoding →

Original recording          Music you hear

- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

11

## Linear Codec

- Suppose $f(x) = Bx^\mathsf{T}$, $B \in R^{m \times d}$.
- Suppose $f'(z) = B'z^\mathsf{T}$, $B' \in R^{d \times m}$.
- We can learn a codec by
- $\min_{B,B'} \sum_{i=1}^{n} \left\|x_i^\mathsf{T} - B'Bx_i^\mathsf{T}\right\|^2$
  - However, there are so many possible candidates $B$ and $B'$!
  - Solving above problem is **hard.**

- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.
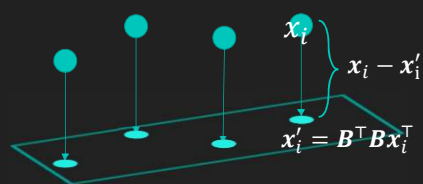
12

## Linear Codec

- We need to put **constraints** on the $B$ and $B'$ to make our problem easier.
- One possible constraint is:
  - $B' = B^\mathsf{T}$
  - $BB' = BB^\mathsf{T} = I$
- Such a codec actually defines an **orthogonal projection of $X$**.
  - Show $B'B$ is an orth. projection matrix

- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.
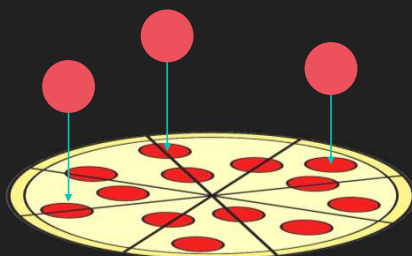
13

## Orthogonal Projection



$x_i$

$x_i - x_i'$

$x_i' = B^\mathsf{T} B x_i^\mathsf{T}$

$z_i = f(x_i) = B x_i^\mathsf{T}$ is called an **embedding** of $x_i$, $B$ is called embedding matrix.

- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

14

## A Pizza Topping Analogy of Embedding



- Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

15

## Minimizing Projection Error

○ $\min\limits_{B,BB^\top=I} \sum_{i=1}^{n} \left|\left|x_i^\top - \textcolor{red}{B^\top B}x_i^\top\right|\right|^2$

  ○ We minimize square error between original data points and its projection.

○ The above problem is equivalent to:

  ○ $\max\limits_{B,BB^\top=I} \mathrm{tr}(BX^\top XB^\top)$

  ○ Live demonstration

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

16

## Minimizing Projection Error

○ $\max\limits_{B,BB^\top=I} \mathrm{tr}(BX^\top XB^\top)$

○ Remarkably, this seemingly complex optimization has an analytical solution:

○ Let $[(\lambda_1, v_1), \dots, (\lambda_m, v_m)]$ be **sorted** eigenvalue and eigenvec of $X^\top X$.

  ○ $\lambda_1 \geq \lambda_2 \dots \geq \lambda_m$

  ○ $\hat{B} = [v_1, v_2, \dots, v_m]^\top$ is an optimal solution, suppose $v_i$ is a **column vector**.

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

17

## Principle Component Analysis

○ As $X$ is a centered dataset,

  ○ $X^\top X = n \cdot \mathrm{cov}[x]$ (PC: show it!)

○ Computing $\hat{B}$ via computing sorted eigenvectors of $\mathrm{cov}[x]$ is called Principle Component Analysis (PCA).

○ Finally, embedding $\hat{f}(x_i) = \hat{B}x_i^\top \in R^m$ is called **PCA embedding** of $x_i$.

  ○ $m$ dimensional "compression" we want!

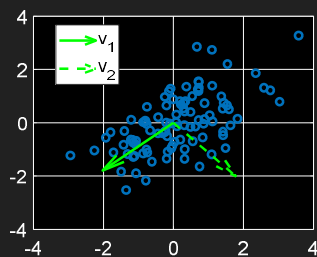• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

18

## Refresh: Eigenvectors and Eigenvalues

- Given a square $n \times n$ matrix $A$, If there exists non-zero vector $v$ such that
- $Av = \lambda v, v \in R^n$
- Then $\lambda$ is an eigenvalue and $v$ is an eigenvector of $A$.

Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.
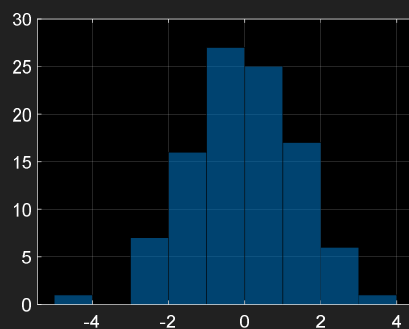
19

## Example, One Cluster



$v_1$ always points at the direction where your dataset has the largest variance!
PC: Intuitively explain why.

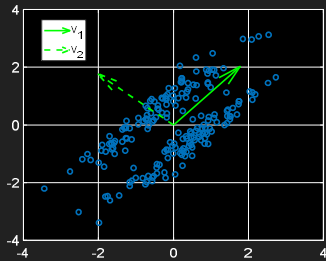Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

20

## Example, Embedding $z = v_1^\mathsf{T} x^\mathsf{T}$



Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.
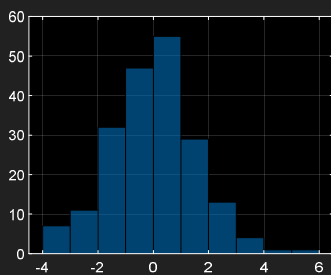
21

## Example, Two Clusters



However, PCA embedding does **not necessarily** preserve clustering information.

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

22

## Example, Embedding $z = v_1^T x^T$



Cluster information **lost** after embedding! Will address this issue in the next lecture.

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

23

## Conclusion

○ Curse of Dimensionality

   ○ $d$ increases, performance may decrease.

○ Principle Component Analysis

   ○ Finding an embedding matrix $\hat{B}$ by computing sorted eigenvalue/vectors of $\mathrm{cov}[x]$.

   ○ PCA Embedding: $\hat{f}(x_i) = \hat{B}x^T$.

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

24