

# Removing Redundancies From Data: Principle Component Analysis

COMS21202, Part III

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

# Objectives

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

# 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  $d > n$

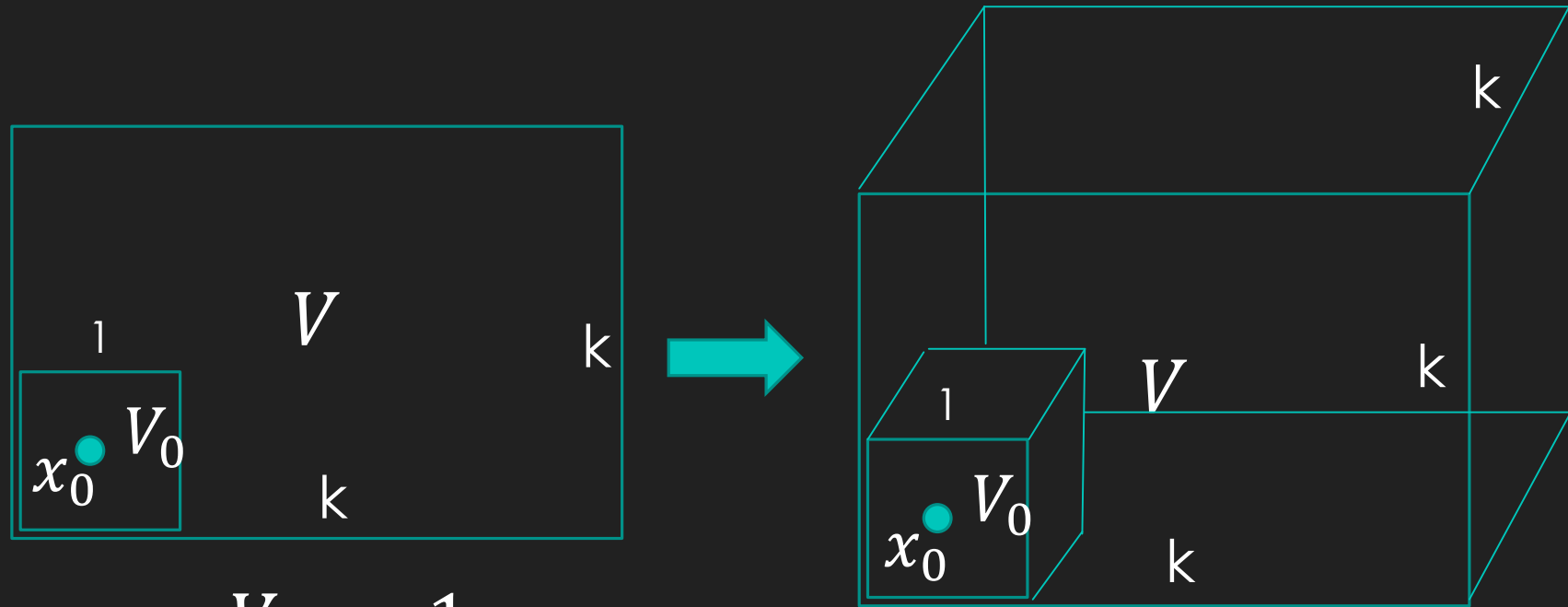
○ ☹ Large  $d$  causes overfitting

○ More ☹ ?

# 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.
- **Reason:** Many geometry concepts can have very different meanings in lower and higher dimensional space.
- One of those is the concept of “locality”.

# The Vanishing Neighborhood



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

$V_0$ : Neighborhood  
volume of  $x_0$

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

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

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

# 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  $x_0$ .

# Reduce the Dimensionality using Feature Transform

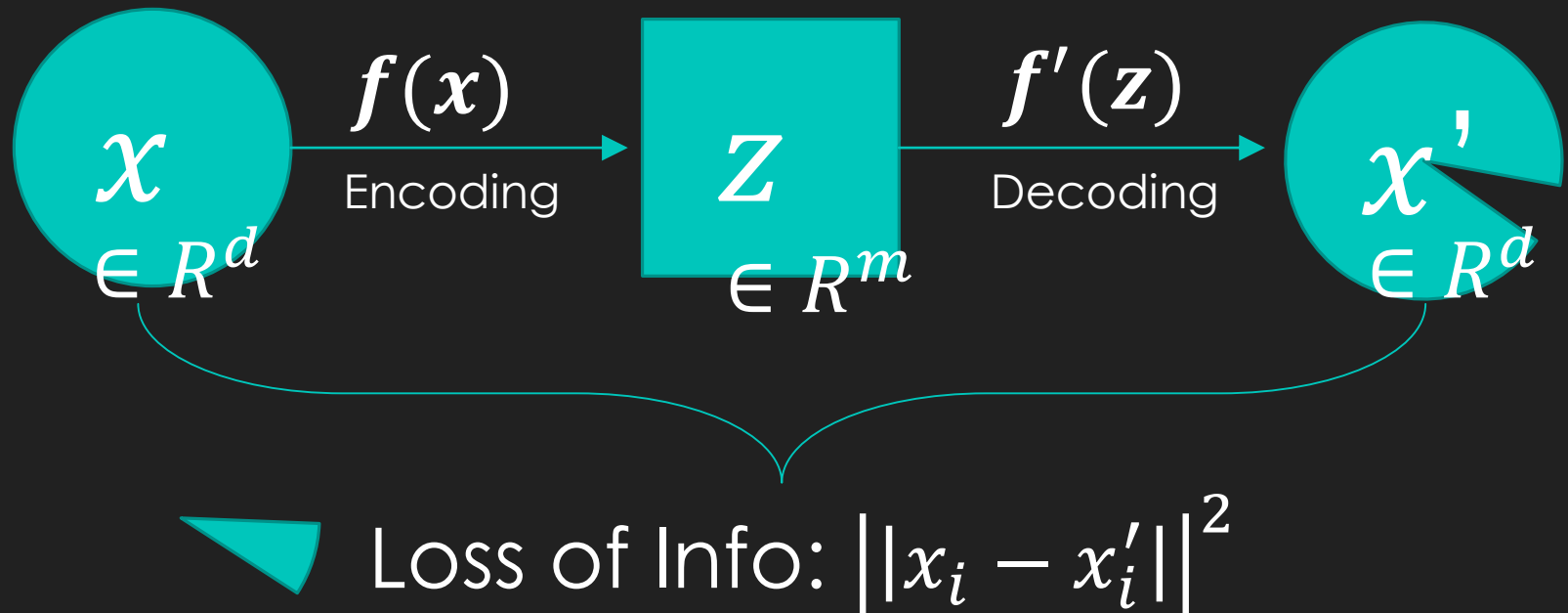
- We want to find a **feature transform**  $f(x) \in R^m$ , where  $m \ll d$ .
  - $f$  transforms original input  $x$  to a subspace as  $R^m \subset R^d$ .
  - We assume our dataset is **centered**:
    - $\frac{1}{n} \sum_{i=1}^n x_i = 0$
    - If dataset  $X'$  is not centered:
      - **Centering**:  $\forall_i x_i = x'_i - \frac{1}{n} \sum_{i=1}^n x'_i$

# Reduce the Dimensionality using Feature Transform

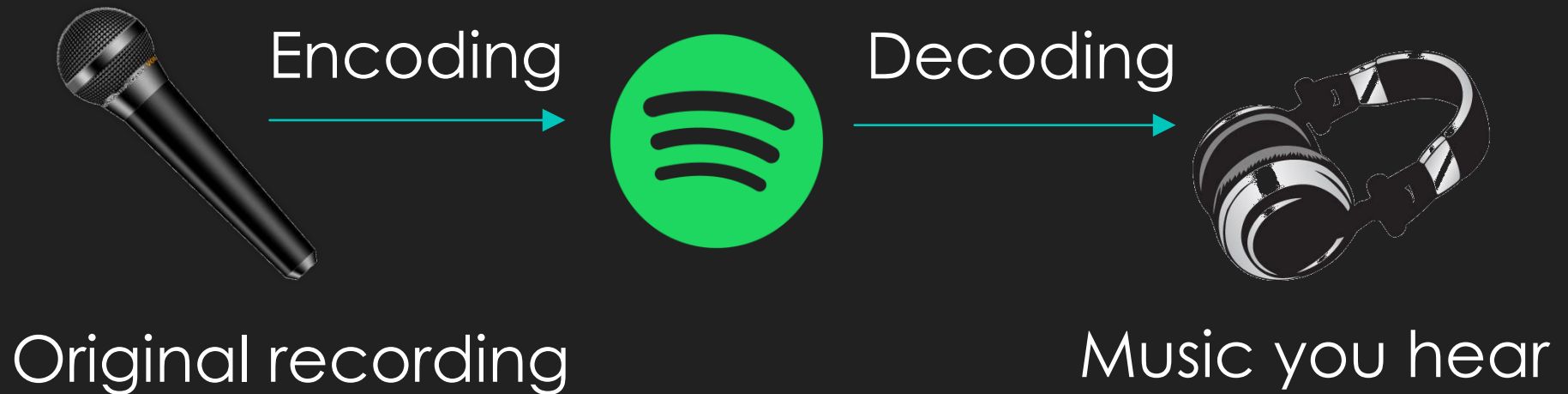
- What is the optimal strategy of selecting  $f(x)$ ?
- Want to reduce dimension using  $f$ .
  - while preserving **as much info as possible!**
- Let's look at this problem from data compression perspective!



# Encoder and Decoder



# Codec



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

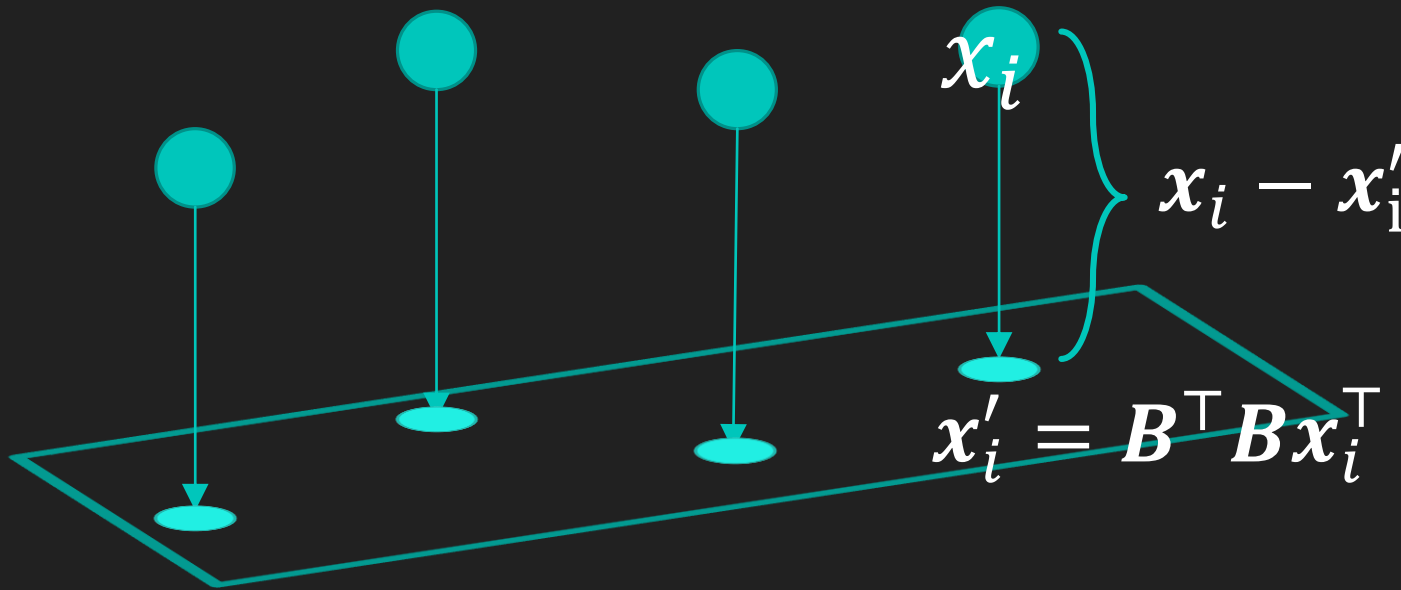
# Linear Codec

- Suppose  $f(\mathbf{x}) = \mathbf{B}\mathbf{x}^\top$ ,  $\mathbf{B} \in \mathbb{R}^{m \times d}$ .
- Suppose  $f'(\mathbf{z}) = \mathbf{B}'\mathbf{z}^\top$ ,  $\mathbf{B}' \in \mathbb{R}^{d \times m}$ .
- We can learn a codec by
- $$\min_{\mathbf{B}, \mathbf{B}'} \sum_{i=1}^n \left\| \mathbf{x}_i^\top - \mathbf{B}'\mathbf{B}\mathbf{x}_i^\top \right\|^2$$
  - However, there are so many possible candidates  $\mathbf{B}$  and  $\mathbf{B}'$ !
  - Solving above problem is **hard**.

# Linear Codec

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

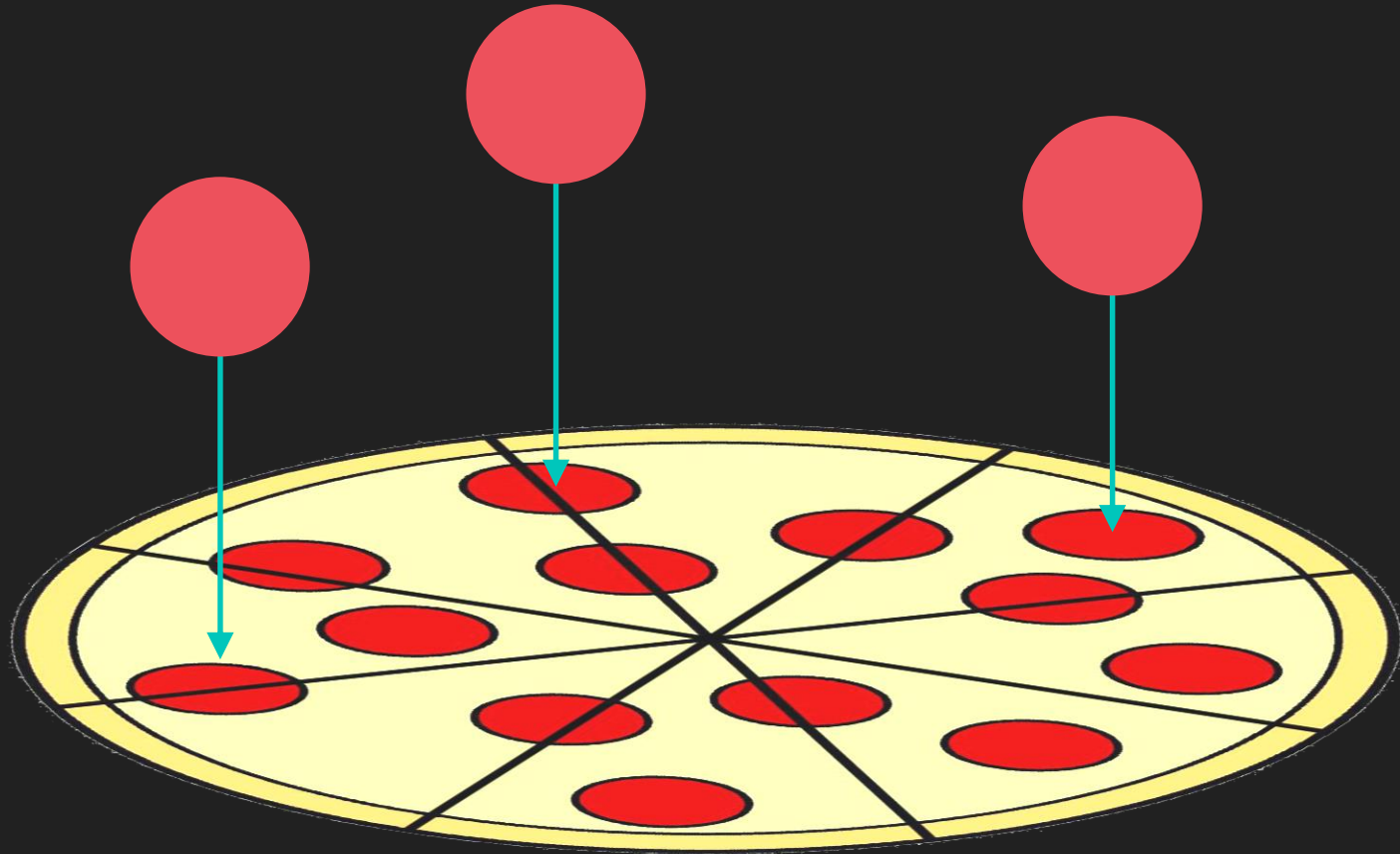
# Orthogonal Projection



$z_i = f(x_i) = Bx_i^T$  is called an **embedding** of  $x_i$ ,  
 $B$  is called embedding matrix.

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

# A Pizza Topping Analogy of Embedding



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

# Minimizing Projection Error

- $\min_{B, BB^T=I} \sum_{i=1}^n \left\| \mathbf{x}_i^T - \mathbf{B}^T \mathbf{B} \mathbf{x}_i^T \right\|^2$

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

- The above problem is equivalent to:

- $\max_{B, BB^T=I} \text{tr}(\mathbf{B} \mathbf{X}^T \mathbf{X} \mathbf{B}^T)$

- Live demonstration

# Minimizing Projection Error

- $\max_{B, BB^T=I} \text{tr}(BX^TXB^T)$
- Remarkably, this seemingly complex optimization has an analytical solution:
- Let  $[(\lambda_1, \mathbf{v}_1), \dots, (\lambda_m, \mathbf{v}_m)]$  be **sorted** eigenvalue and eigenvec of  $X^TX$ .
  - $\lambda_1 \geq \lambda_2 \dots \geq \lambda_m$
  - $\hat{B} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m]^T$  is **an optimal solution**, suppose  $\mathbf{v}_i$  is a **column vector**.



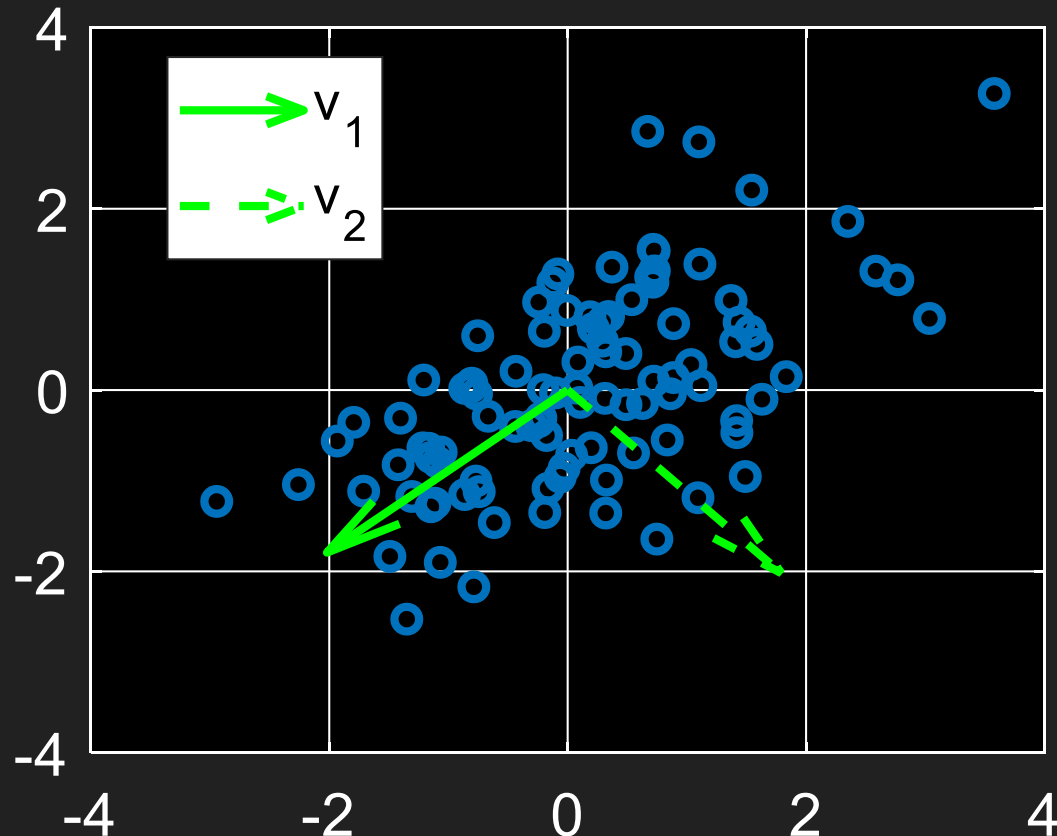
# Principle Component Analysis

- As  $\mathbf{X}$  is a centered dataset,
  - $\mathbf{X}^T \mathbf{X} = n \cdot \text{cov}[\mathbf{x}]$  (PC: show it!)
- Computing  $\hat{\mathbf{B}}$  via computing sorted eigenvectors of  $\text{cov}[\mathbf{x}]$  is called Principle Component Analysis (PCA).
- Finally, embedding  $\hat{\mathbf{f}}(\mathbf{x}_i) = \hat{\mathbf{B}} \mathbf{x}_i^T \in \mathbb{R}^m$  is called **PCA embedding** of  $\mathbf{x}_i$ .
  - $m$  dimensional “compression” we want!

# 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  $A$ .

# Example, One Cluster

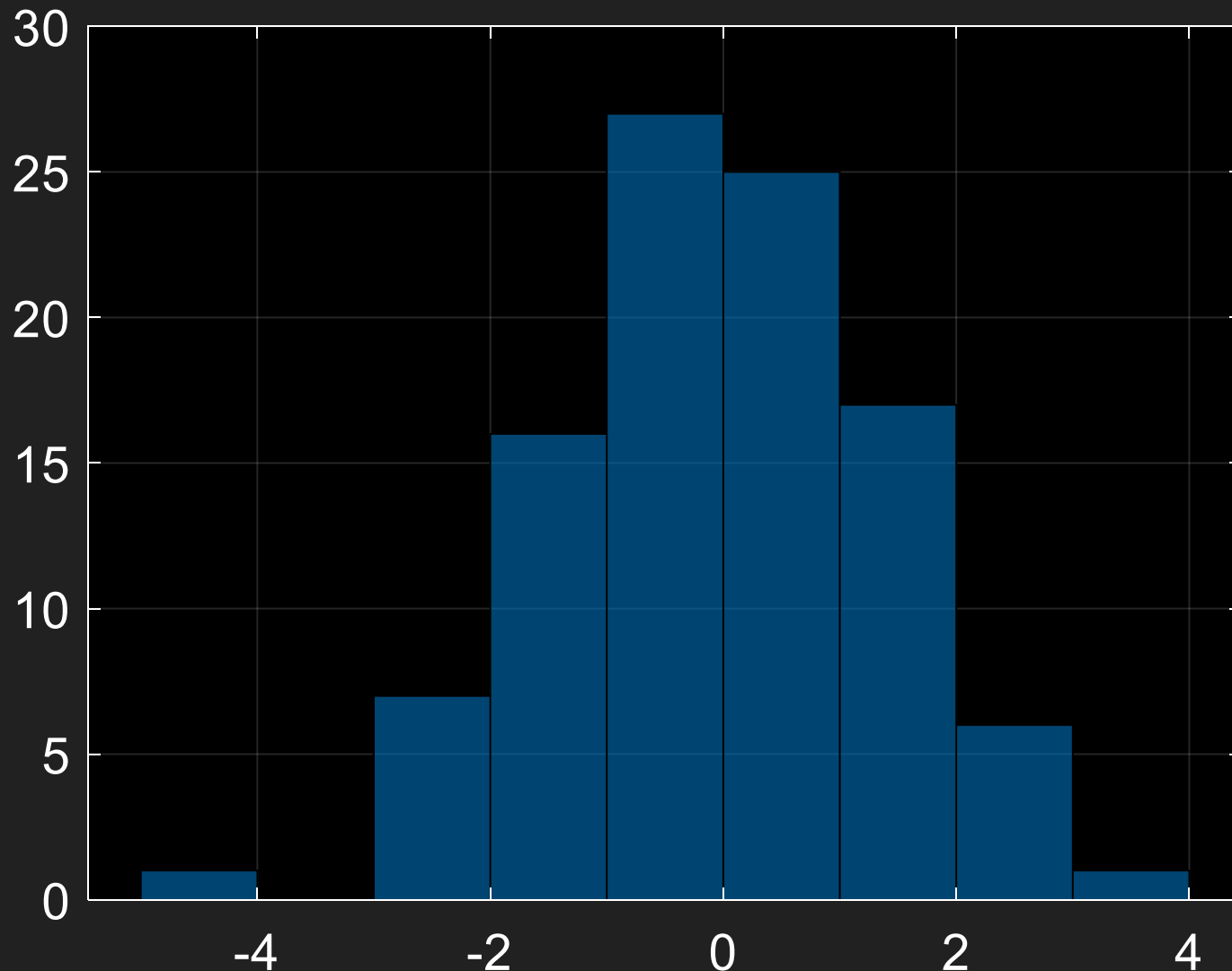


$v_1$  always points at the direction where your dataset has the largest variance!

PC: Intuitively explain why.

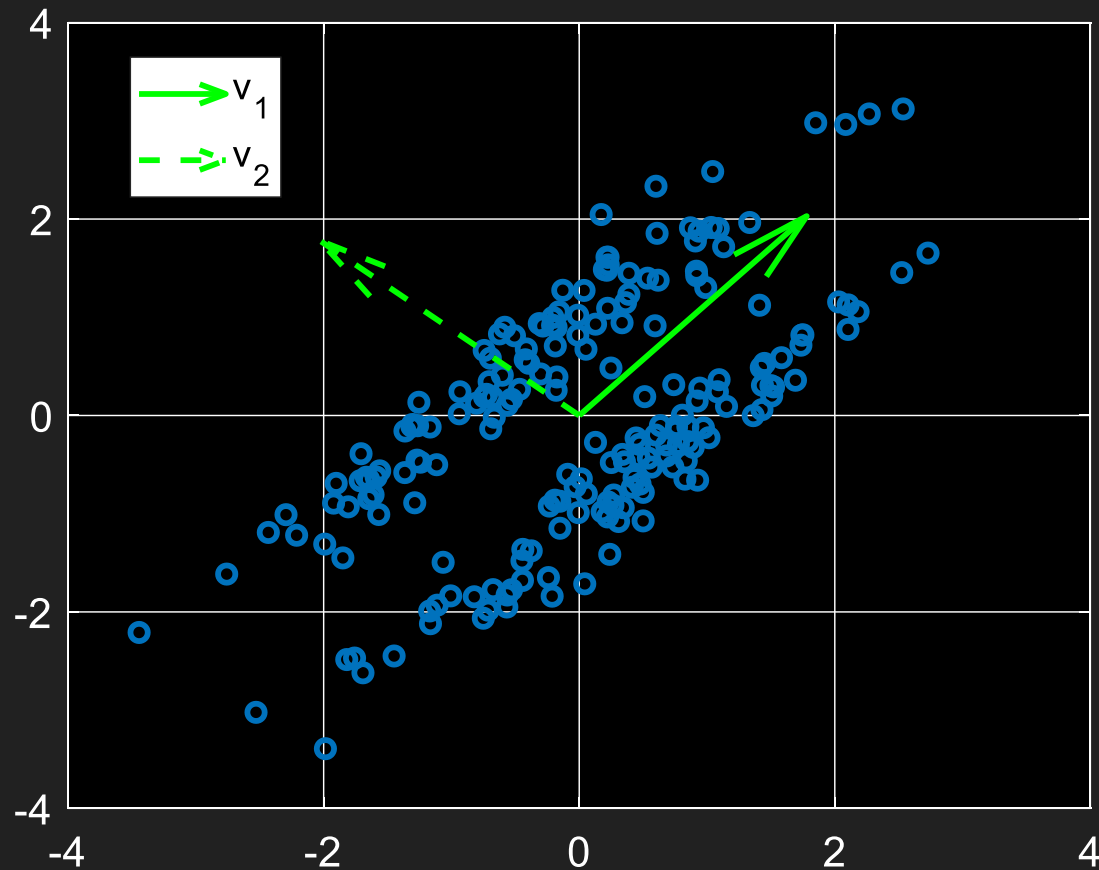
- Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

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



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

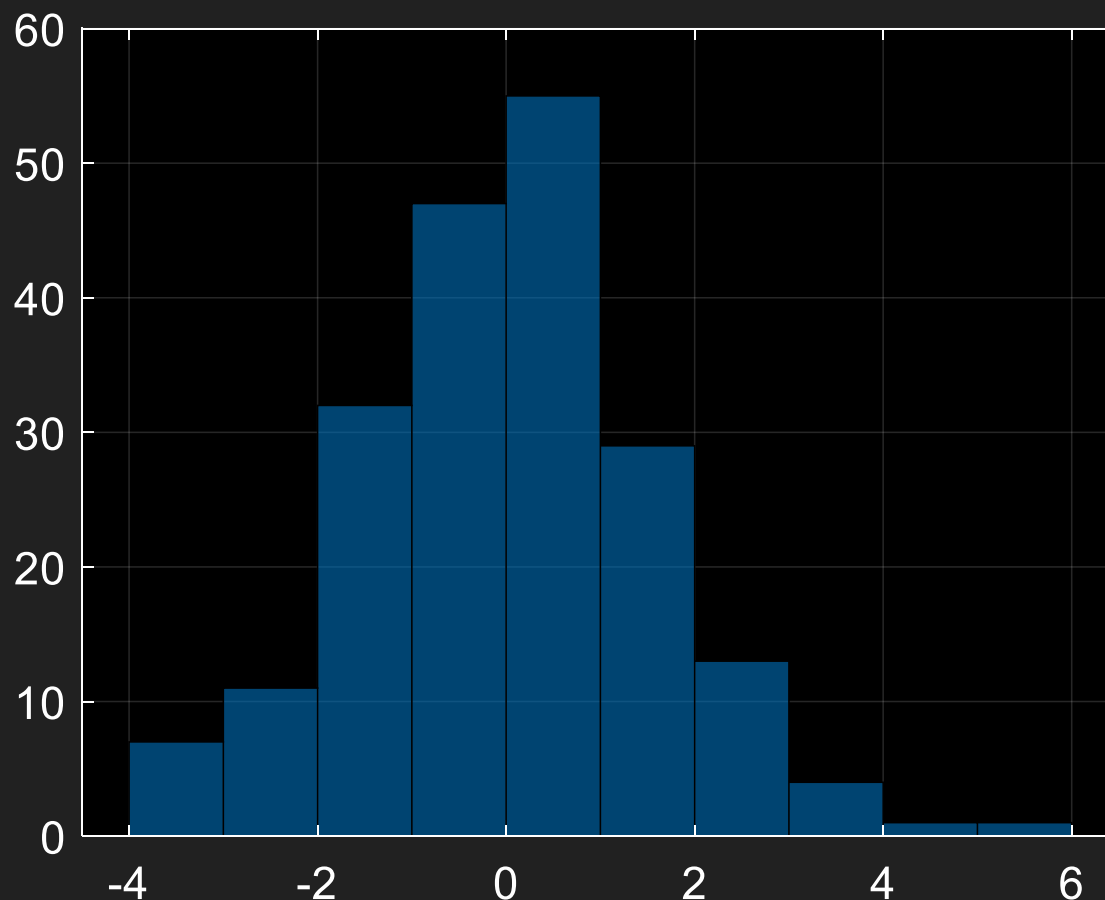
# Example, Two Clusters



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

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

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



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

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

# Conclusion

## ○ Curse of Dimensionality

- $d$  increases, performance may decrease.

## ○ Principle Component Analysis

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

- PCA Embedding:  $\hat{\mathbf{f}}(\mathbf{x}_i) = \hat{\mathbf{B}}\mathbf{x}^\top$ .