# Compactifying Fish

## Principal Component Analysis and Dimension Reduction

Eli Griffiths

# Simple Data is Not So Simple

# Simple Data is Not So Simple
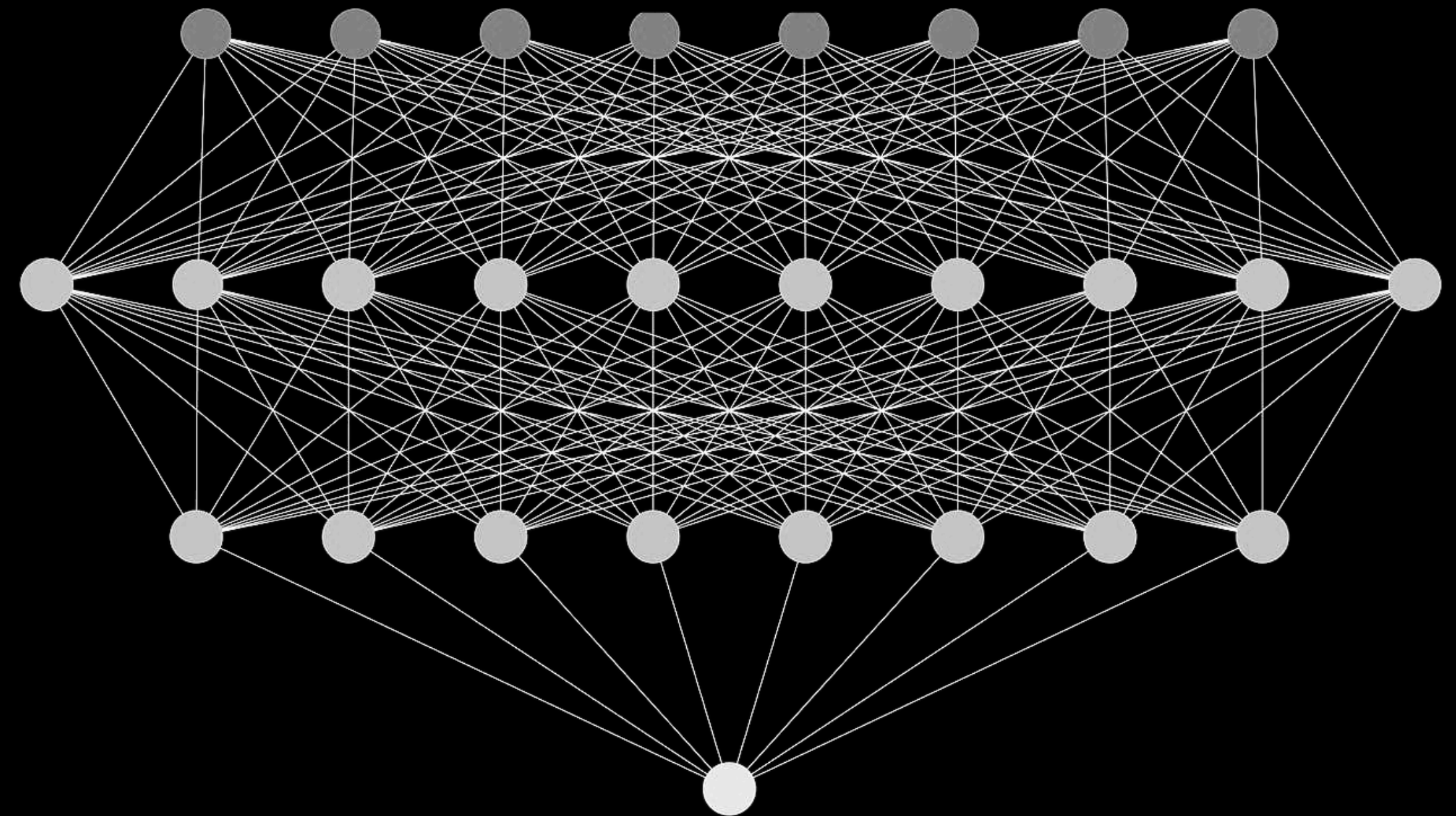


32

32

# Simple Data is Not So Simple



$32$

$32$

$$32 \times 32 \times (\textcolor{red}{1} + \textcolor{blue}{1} + \textcolor{green}{1}) = 3,072$$
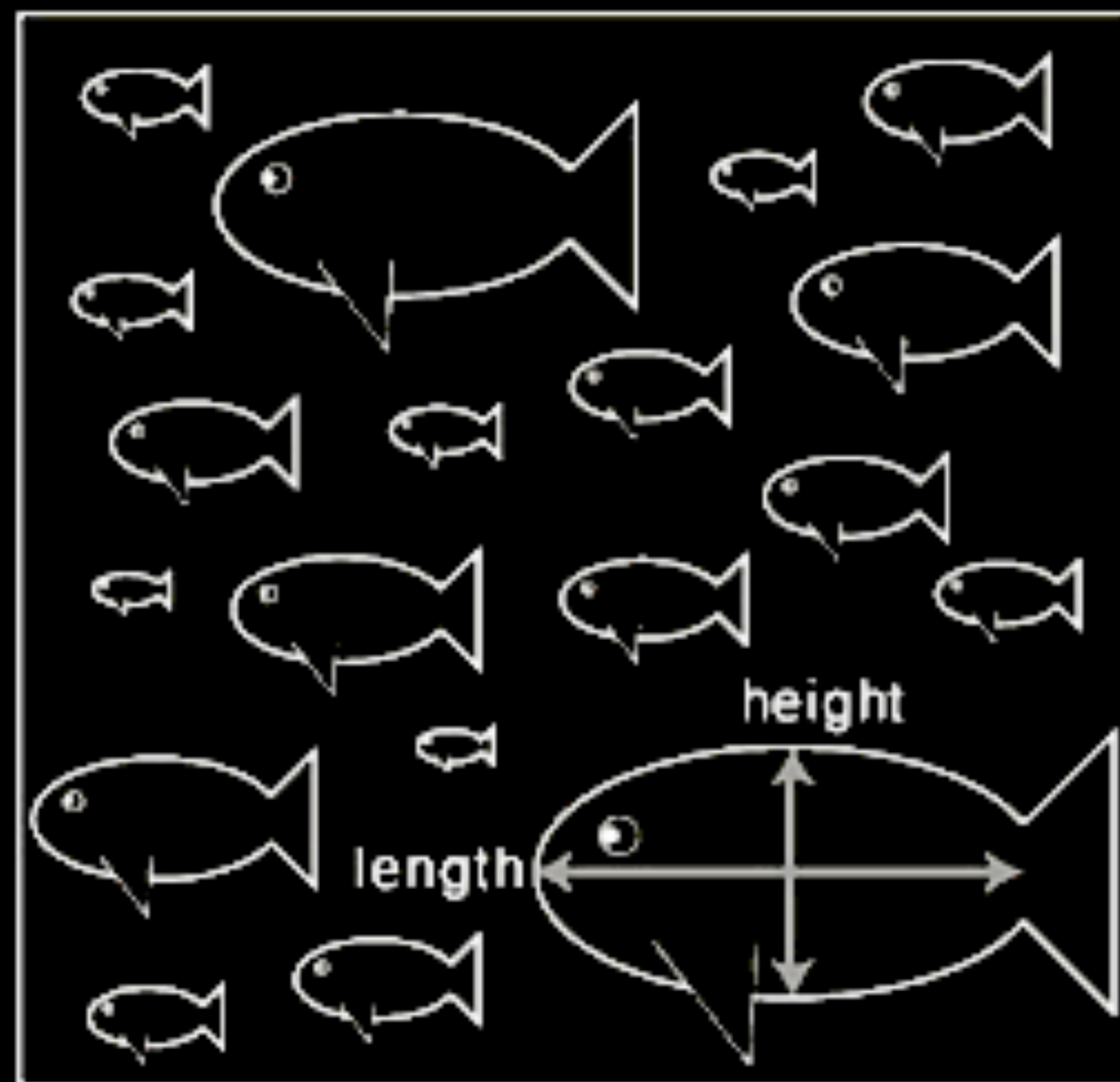
# Why Dimension Reduction?

- **Data Science** – Analysis of high featured datasets

- **Machine Learning** – Dataset simplification

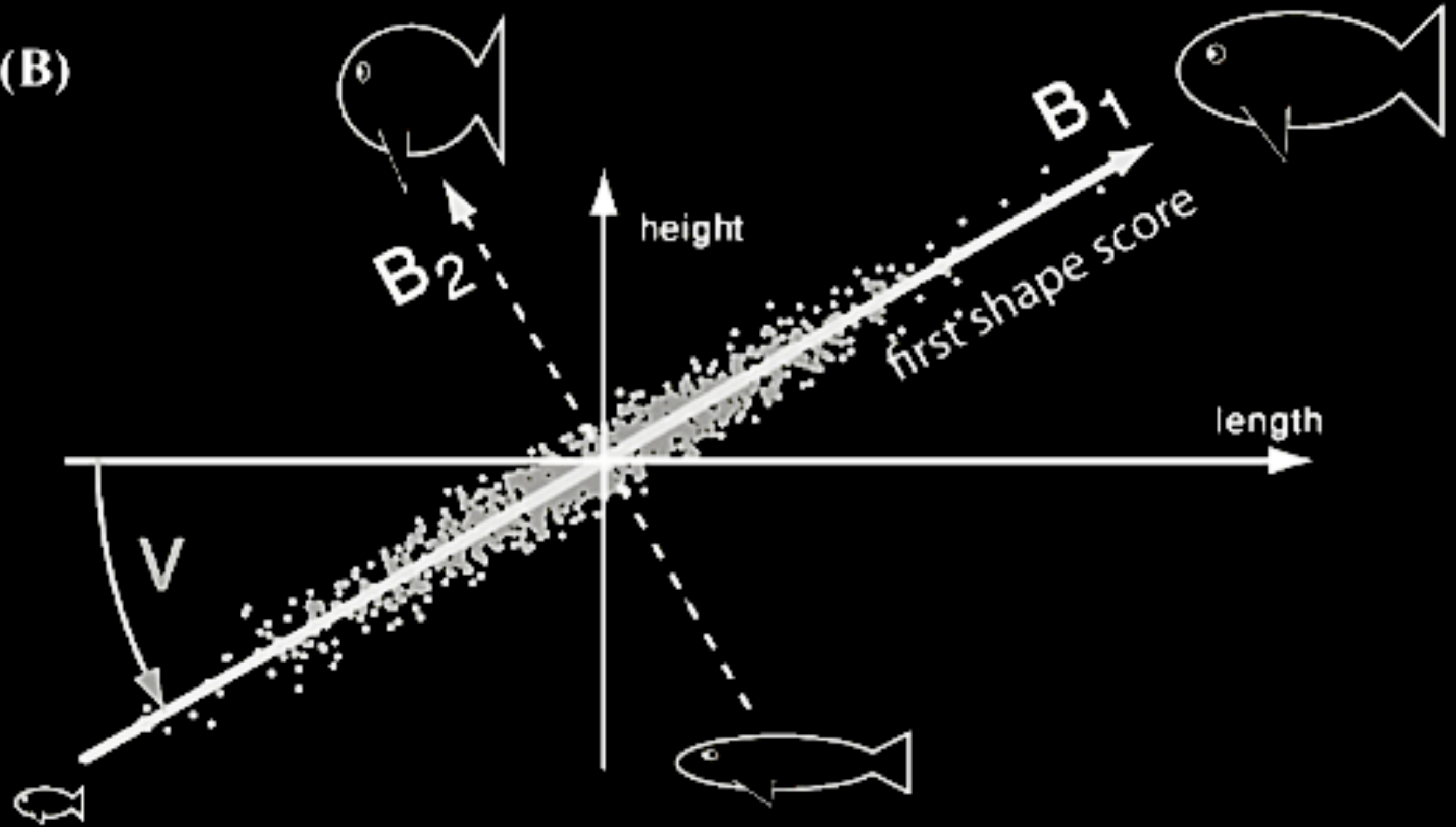- **Neuroscience** – Neuron potentials and activation
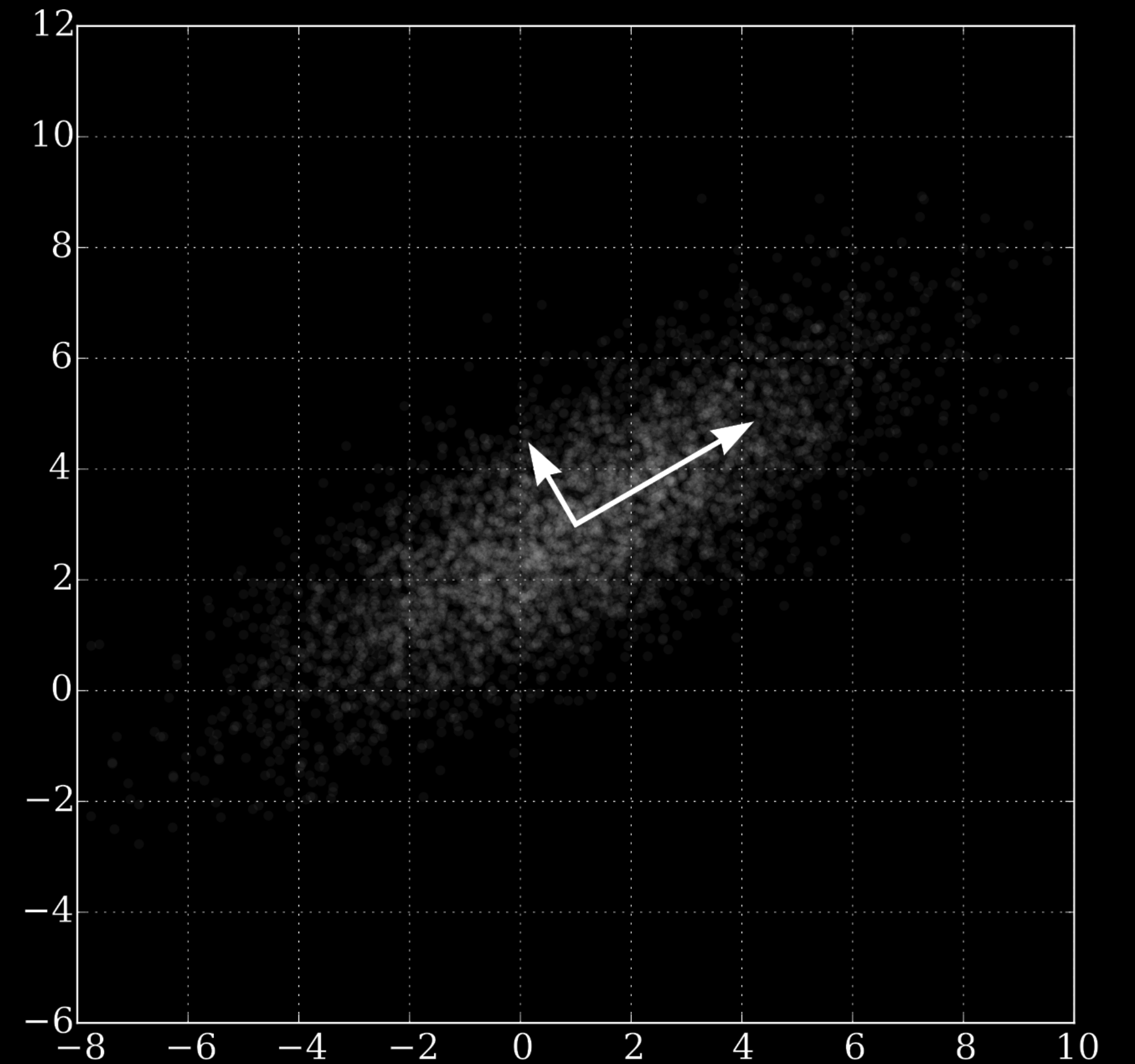
# Reduction via Correlation

# Principal Components

- Identify a set of correlations

- Pick the strongest ones to build *axes*

- Project the data onto these axes

# Capturing Correlation

$$\mathbf{X} = \begin{bmatrix} X_1 & X_2 & \cdots & X_n \end{bmatrix}^T$$

# Capturing Correlation

$$\mathbf{X} = \begin{bmatrix} X_1 & X_2 & \cdots & X_n \end{bmatrix}^T$$

$$\mathbf{K}_{ij} = \mathrm{Cov}[X_i, X_j] = \sigma_{ij}$$

# Capturing Correlation

$$\mathbf{X} = \begin{bmatrix} X_1 & X_2 & \cdots & X_n \end{bmatrix}^T$$

$$\mathbf{K}_{ij} = \text{Cov}[X_i, X_j] = \sigma_{ij}$$

$$\mathbf{K} = \begin{bmatrix} \sigma_{1,1} & \sigma_{1,2} & \cdots & \sigma_{1,n} \\ \sigma_{2,1} & \sigma_{2,2} & & \sigma_{2,n} \\ \vdots & & \ddots & \vdots \\ \sigma_{n,1} & \sigma_{n,2} & \cdots & \sigma_{n,n} \end{bmatrix}$$

# Symmetric Matrix

# Symmetric Matrix

$$\mathrm{Cov}[X_i, X_j] = \mathrm{Cov}[X_j, X_i]$$

# Symmetric Matrix

$$\text{Cov}[X_i, X_j] = \text{Cov}[X_j, X_i]$$

$$\mathbf{K} = \begin{bmatrix} \sigma_{1,1} & \sigma_{1,2} & \cdots & \sigma_{1,n} \\ \sigma_{1,2} & \sigma_{2,2} & & \sigma_{2,n} \\ \vdots & & \ddots & \vdots \\ \sigma_{1,n} & \sigma_{2,n} & \cdots & \sigma_{n,n} \end{bmatrix}$$

# Spectral Theorem

If $A$ is symmetric, there exists an orthonormal ***basis*** of eigenvectors of $A$

# The Projection

$$\mathbf{K}\, x_i = \lambda_i x_i \qquad \lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_d \geq \ldots \geq \lambda_n$$

# The Projection

$$\mathbf{K}\, x_i = \lambda_i x_i \qquad \lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_d \geq \ldots \geq \lambda_n$$

$$\mathbf{P} = \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \cdots & x_d \\ | & | & & | \end{bmatrix}$$

# A Concrete Example

- MNIST Handwriting Dataset

- Comprised of 28 by 28 grayscale images

- Has 784 "features"

Original          $d = 50$          $d = 100$          $d = 300$          $d = 500$

Thank You