

# 10-301/601: Introduction to Machine Learning

## Lecture 24: Dimensionality Reduction

Henry Chai

6/5/25

# Front Matter

- Announcements
  - HW5 released on 6/3, due 6/6 (tomorrow) at 11:59 PM
  - HW6 to be released on 6/6 (tomorrow), due 6/10 at 11:59 PM
  - Quiz 3 on 6/6 (tomorrow) at 11:00 AM in BH A36

# Learning Paradigms

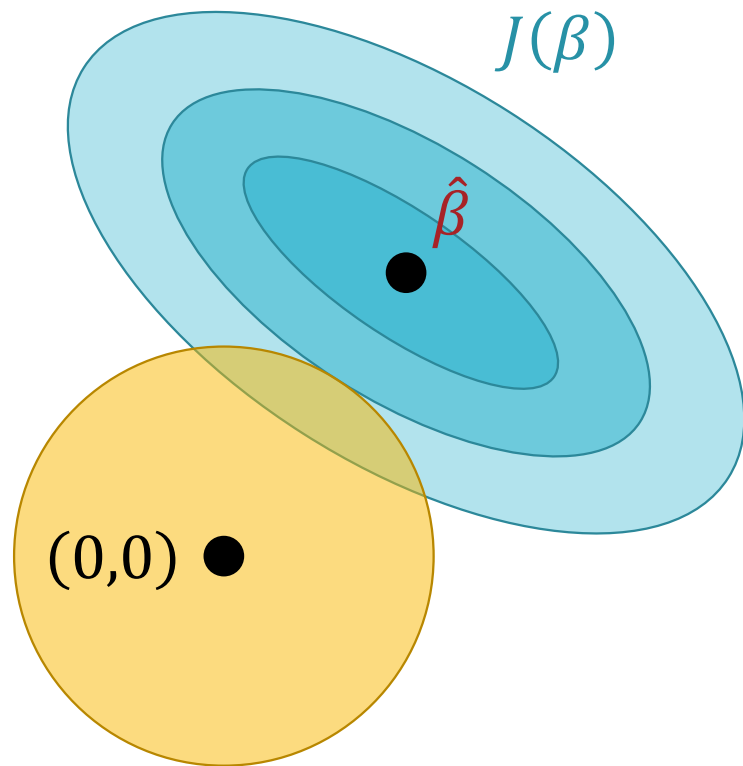
- Supervised learning -  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$ 
  - Regression -  $y^{(n)} \in \mathbb{R}$
  - Classification -  $y^{(n)} \in \{1, \dots, C\}$
- Unsupervised learning -  $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$ 
  - Clustering
  - Dimensionality reduction

# Learning Paradigms

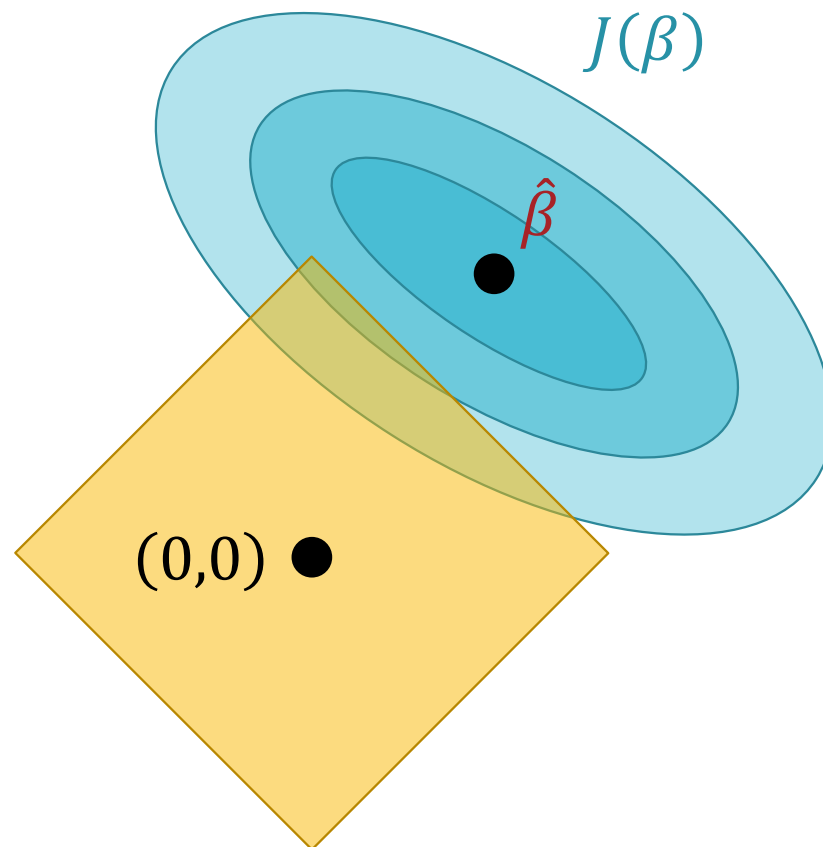
- Supervised learning -  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$ 
  - Regression -  $y^{(n)} \in \mathbb{R}$
  - Classification -  $y^{(n)} \in \{1, \dots, C\}$
- Unsupervised learning -  $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$ 
  - Clustering
  - **Dimensionality reduction**

# Dimensionality Reduction

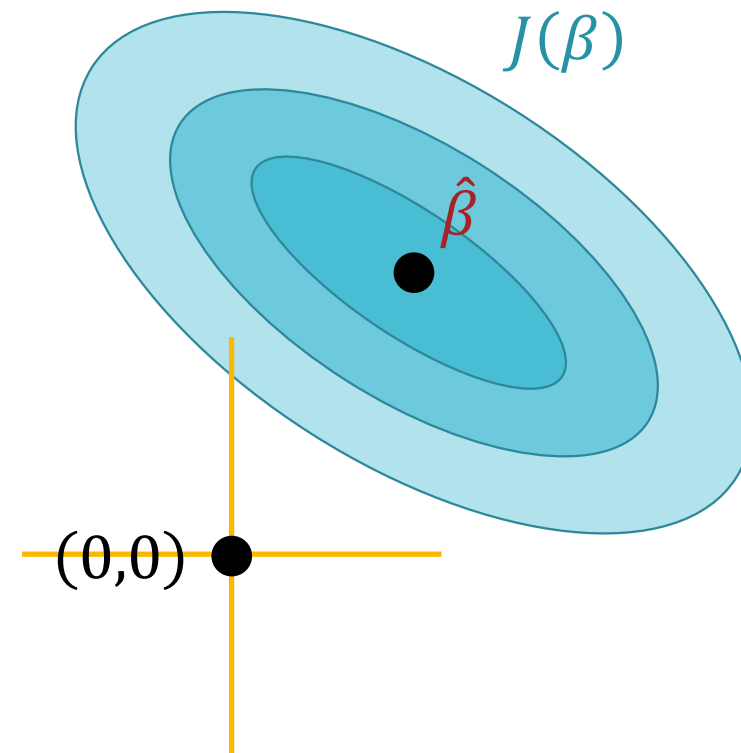
- Goal: given some unlabeled data set, learn a latent (typically lower-dimensional) representation
- Use cases:
  - Reducing computational cost (runtime, storage, etc...)
  - Improving generalization
  - Visualizing data
- Applications:
  - High-resolution images/videos
  - Text data
  - Financial or transaction data



Ridge or  $L_2$

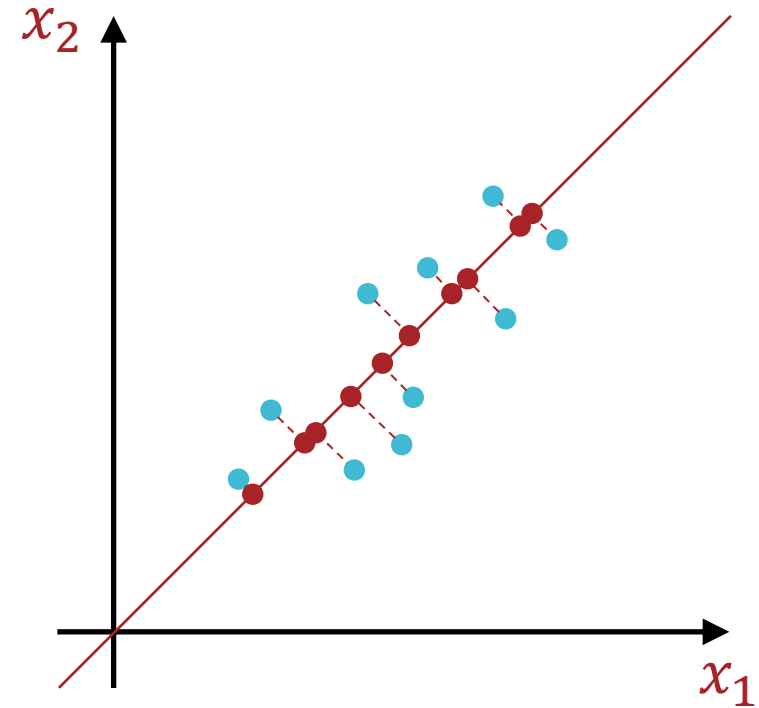
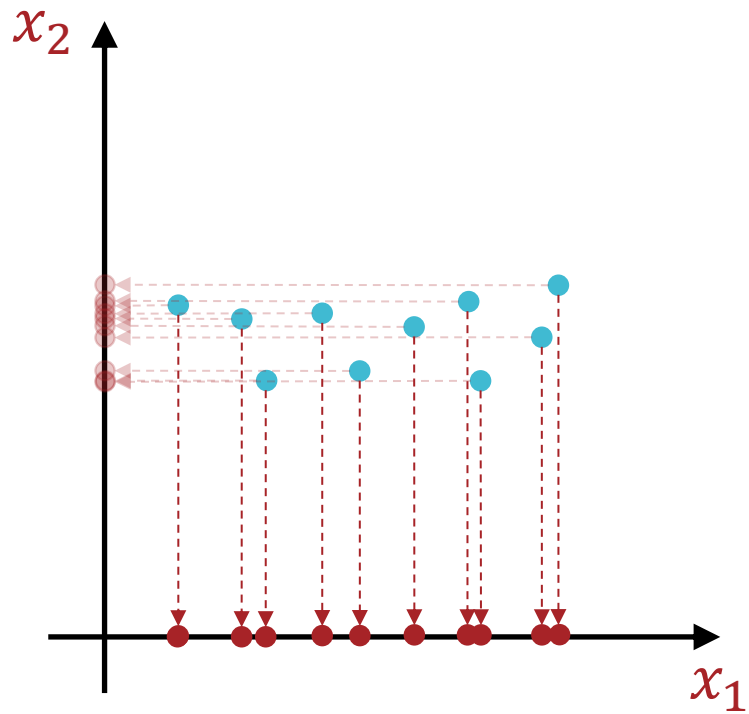


Lasso or  $L_1$

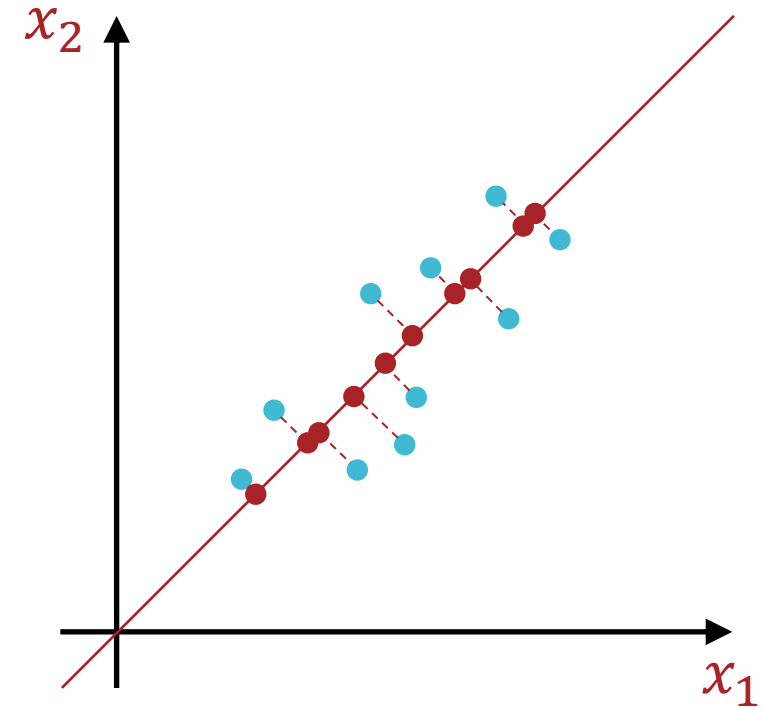
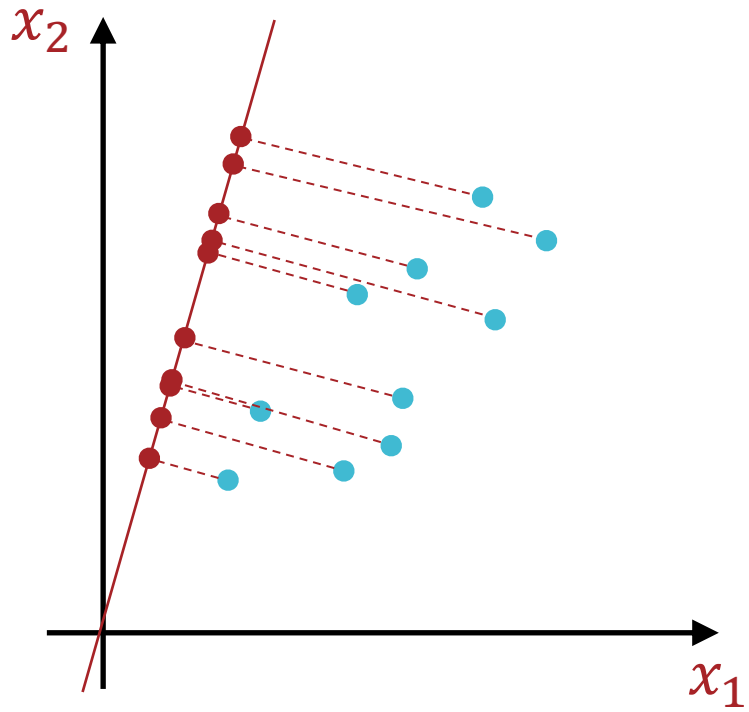


$L_0$

Recall:  $L_1$  (or  $L_0$ ) Regularization

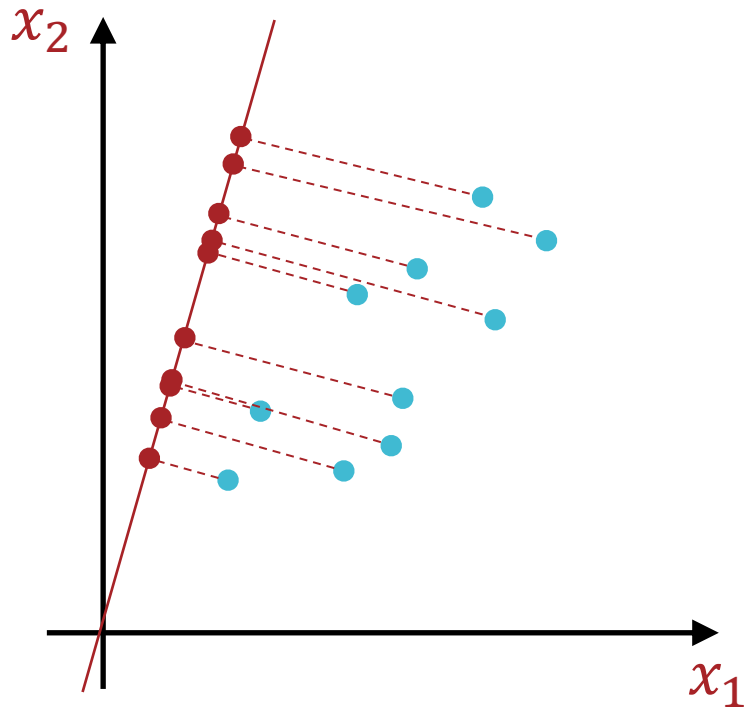


# Feature Elimination

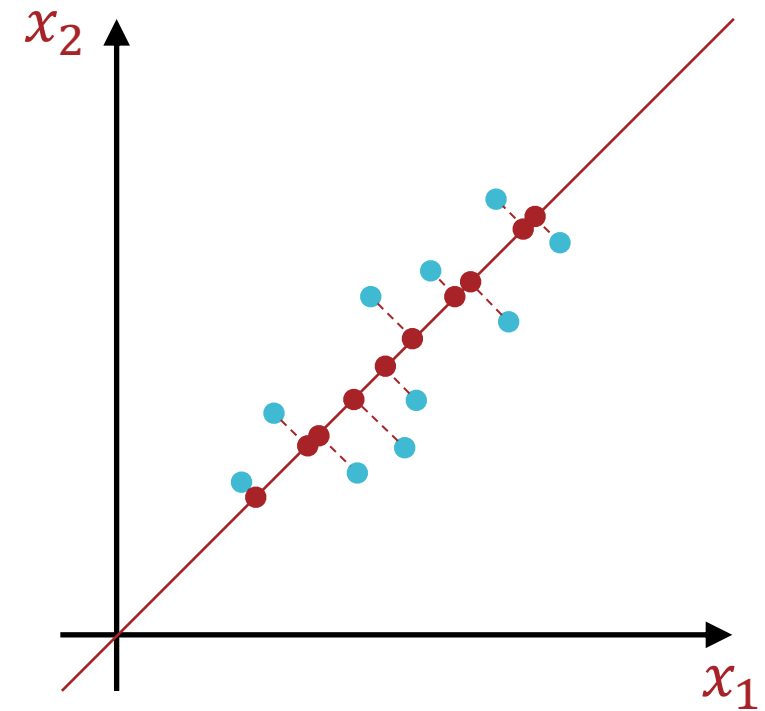


# Feature Reduction





Option A



Option B

Which projection do you prefer?

# Centering the Data

- To be consistent, we will constrain principal components to be *orthogonal unit vectors* that begin at the origin
- Preprocess data to be centered around the origin:

$$1. \boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)}$$

$$2. \tilde{\mathbf{x}}^{(n)} = \mathbf{x}^{(n)} - \boldsymbol{\mu} \forall n$$

$$3. X = \begin{bmatrix} \tilde{\mathbf{x}}^{(1)T} \\ \tilde{\mathbf{x}}^{(2)T} \\ \vdots \\ \tilde{\mathbf{x}}^{(N)T} \end{bmatrix}$$

# Reconstruction Error

- The projection of  $\tilde{\mathbf{x}}^{(n)}$  onto a vector  $\mathbf{v}$  is

$$\mathbf{z}^{(n)} = \left( \frac{\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}}{\|\mathbf{v}\|_2} \right) \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$$

Length of projection

Direction of projection

# Reconstruction Error

- The projection of  $\tilde{\mathbf{x}}^{(n)}$  onto a unit vector  $\mathbf{v}$  is

$$\mathbf{z}^{(n)} = (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}$$

$$\hat{\mathbf{v}} = \operatorname{argmin}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \sum_{n=1}^N \|\tilde{\mathbf{x}}^{(n)} - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}\|_2^2$$

$$\|\tilde{\mathbf{x}}^{(n)} - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}\|_2^2$$

$$= \tilde{\mathbf{x}}^{(n)T} \tilde{\mathbf{x}}^{(n)} - 2(\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}^T \tilde{\mathbf{x}}^{(n)} + (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}^T \mathbf{v}$$

$$= \tilde{\mathbf{x}}^{(n)T} \tilde{\mathbf{x}}^{(n)} - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}^T \tilde{\mathbf{x}}^{(n)}$$

$$= \|\tilde{\mathbf{x}}^{(n)}\|_2^2 - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)})^2$$

Minimizing the  
Reconstruction  
Error



Maximizing the  
Variance

$$\hat{\mathbf{v}} = \operatorname{argmin}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \sum_{n=1}^N \|\tilde{\mathbf{x}}^{(n)} - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}\|_2^2$$

$$= \operatorname{argmin}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \sum_{n=1}^N \|\tilde{\mathbf{x}}^{(n)}\|_2^2 - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)})^2$$

$$= \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \sum_{n=1}^N (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)})^2 \quad \leftarrow \text{Variance of projections } (\tilde{\mathbf{x}}^{(n)} \text{ are centered})$$

$$= \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \mathbf{v}^T \left( \sum_{n=1}^N \tilde{\mathbf{x}}^{(n)} \tilde{\mathbf{x}}^{(n)T} \right) \mathbf{v}$$

$$= \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \mathbf{v}^T (X^T X) \mathbf{v}$$

# Maximizing the Variance

$$\hat{\mathbf{v}} = \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2^2 = 1} \mathbf{v}^T (X^T X) \mathbf{v}$$

$$\begin{aligned} \mathcal{L}(\mathbf{v}, \lambda) &= \mathbf{v}^T (X^T X) \mathbf{v} - \lambda (\|\mathbf{v}\|_2^2 - 1) \\ &= \mathbf{v}^T (X^T X) \mathbf{v} - \lambda (\mathbf{v}^T \mathbf{v} - 1) \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}} = (X^T X) \mathbf{v} - \lambda \mathbf{v}$$

$$\rightarrow (X^T X) \hat{\mathbf{v}} - \lambda \hat{\mathbf{v}} = 0 \rightarrow (X^T X) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}}$$

- $\hat{\mathbf{v}}$  is an eigenvector of  $X^T X$  and  $\lambda$  is the corresponding eigenvalue! But which one?

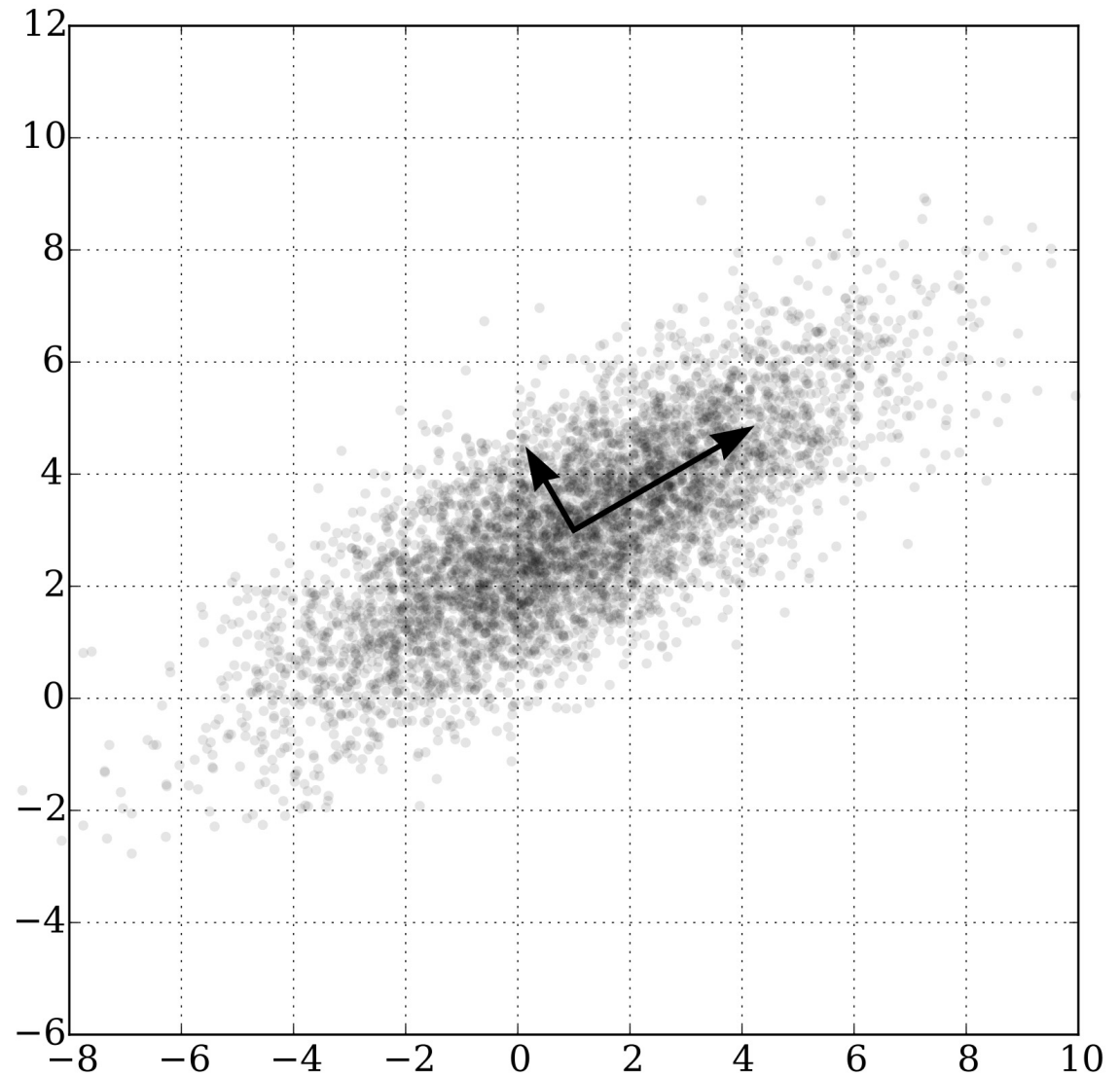
# Maximizing the Variance

$$\hat{\mathbf{v}} = \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \mathbf{v}^T (X^T X) \mathbf{v}$$

$$(X^T X) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}} \rightarrow \hat{\mathbf{v}}^T (X^T X) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}}^T \hat{\mathbf{v}} = \lambda$$

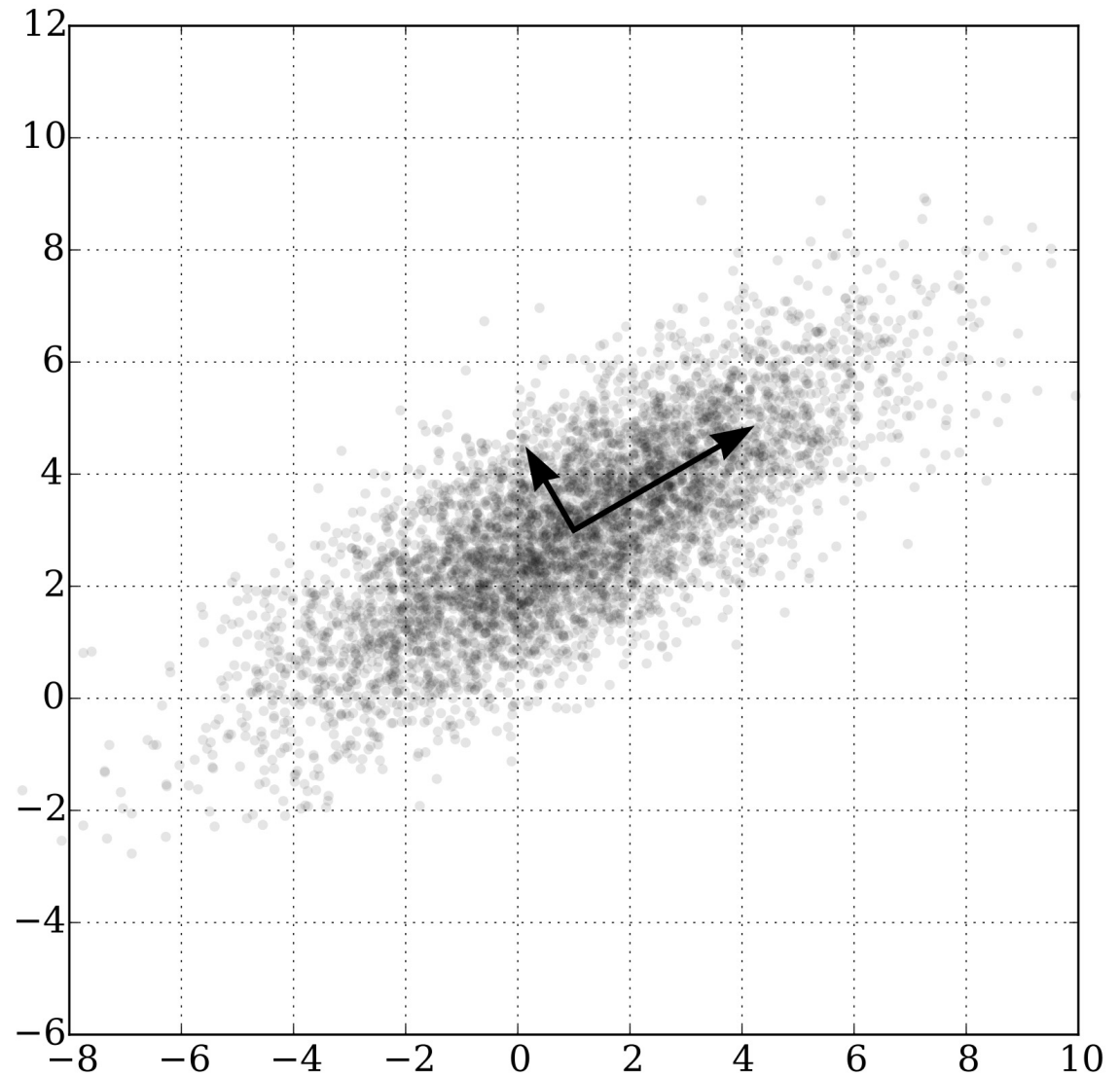
- The first principal component is the eigenvector  $\hat{\mathbf{v}}_1$  that corresponds to the largest eigenvalue  $\lambda_1$
- The second principal component is the eigenvector  $\hat{\mathbf{v}}_2$  that corresponds to the second largest eigenvalue  $\lambda_2$ 
  - $\hat{\mathbf{v}}_1$  and  $\hat{\mathbf{v}}_2$  are orthogonal
- Etc ...
- $\lambda_i$  is a measure of how much variance falls along  $\hat{\mathbf{v}}_i$

# Principal Components: Example





How can we  
efficiently find  
principal  
components  
(eigenvectors)?



# Singular Value Decomposition (SVD) for PCA

- Every real-valued matrix  $X \in \mathbb{R}^{N \times D}$  can be expressed as

$$X = USV^T$$

where:

1.  $U \in \mathbb{R}^{N \times N}$  - columns of  $U$  are eigenvectors of  $XX^T$
2.  $V \in \mathbb{R}^{D \times D}$  - columns of  $V$  are eigenvectors of  $X^T X$
3.  $S \in \mathbb{R}^{N \times D}$  - diagonal matrix whose entries are the eigenvalues of  $X \rightarrow$  squared entries are the eigenvalues of  $XX^T$  and  $X^T X$

# PCA Algorithm

- Input:  $\mathcal{D} = \{(\mathbf{x}^{(n)})\}_{n=1}^N, \rho$ 
  1. Center the data
  2. Use SVD to compute the eigenvalues and eigenvectors of  $\mathbf{X}^T \mathbf{X}$
  3. Collect the top  $\rho$  eigenvectors (corresponding to the  $\rho$  largest eigenvalues),  $\mathbf{V}_\rho \in \mathbb{R}^{D \times \rho}$
  4. Project the data into the space defined by  $\mathbf{V}_\rho$ ,  $\mathbf{Z} = \mathbf{X} \mathbf{V}_\rho$
- Output:  $\mathbf{Z}$ , the transformed (potentially lower-dimensional) data

## How many PCs should we use?

- Input:  $\mathcal{D} = \{(\mathbf{x}^{(n)})\}_{n=1}^N, \rho$ 
  1. Center the data
  2. Use SVD to compute the eigenvalues and eigenvectors of  $\mathbf{X}^T \mathbf{X}$
  3. Collect the top  $\rho$  eigenvectors (corresponding to the  $\rho$  largest eigenvalues),  $\mathbf{V}_\rho \in \mathbb{R}^{D \times \rho}$
  4. Project the data into the space defined by  $\mathbf{V}_\rho$ ,  $\mathbf{Z} = \mathbf{X}\mathbf{V}_\rho$
- Output:  $\mathbf{Z}$ , the transformed (potentially lower-dimensional) data

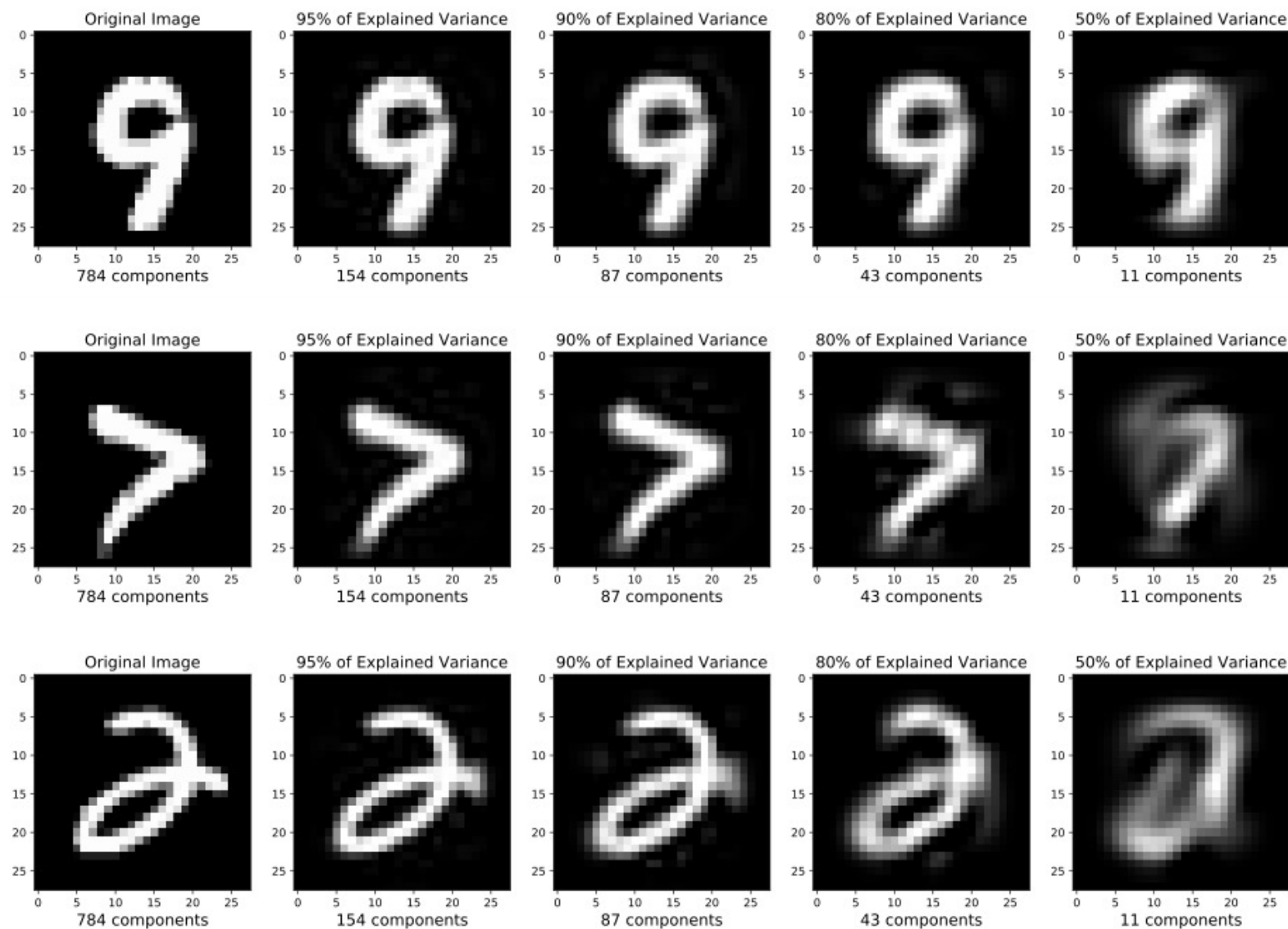
## Choosing the number of PCs

- Define a percentage of explained variance for the  $i^{\text{th}}$  PC:

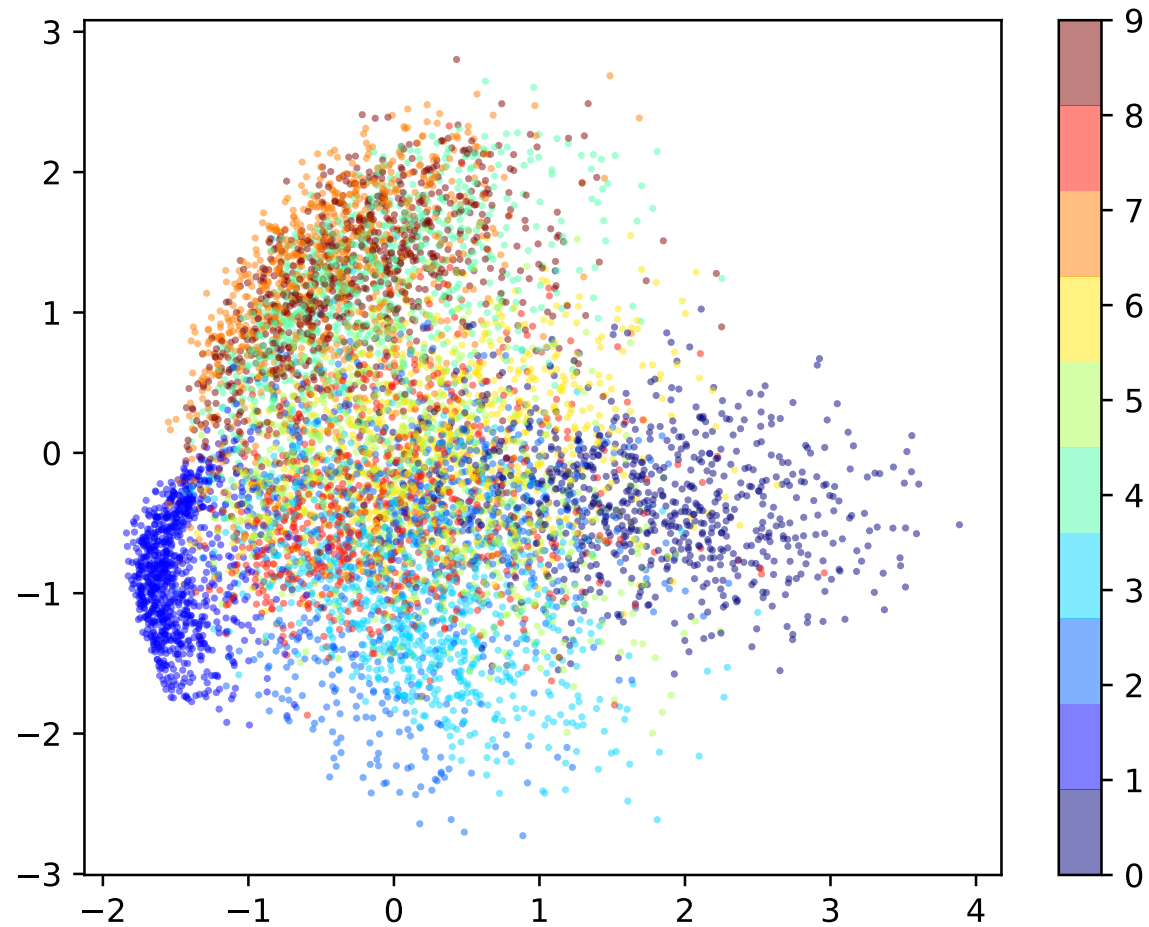
$$\lambda_i / \sum \lambda_j$$

- Select all PCs above some threshold of explained variance, e.g., 5%
- Keep selecting PCs until the total explained variance exceeds some threshold, e.g., 90%
- Evaluate on some downstream metric

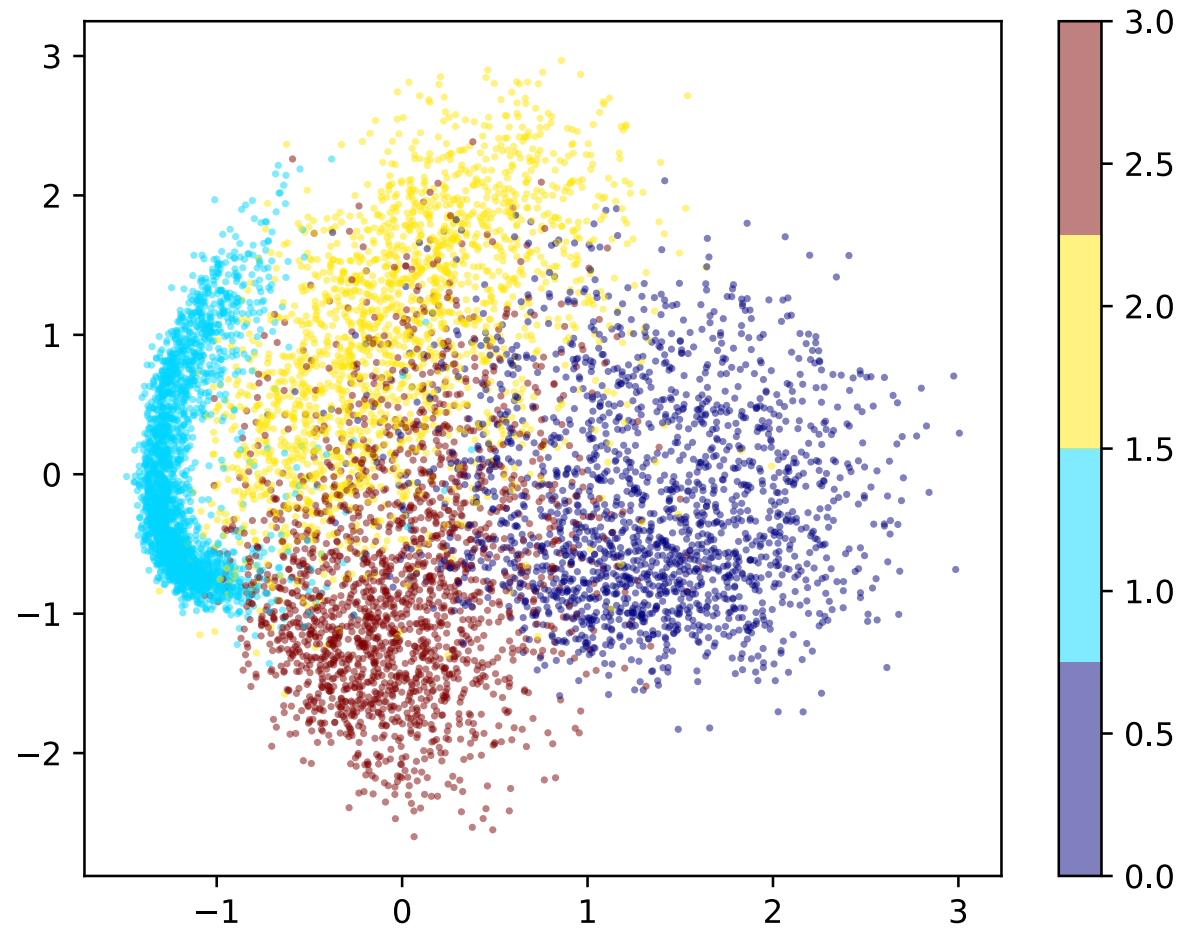
# PCA Example: MNIST Digits



# PCA Example: MNIST Digits



# PCA Example: MNIST Digits





# Key Takeaways

- PCA finds an orthonormal basis where the first principal component maximizes the variance  $\Leftrightarrow$  minimizes the reconstruction error
- Autoencoders use neural networks to automatically learn a latent representation that minimizes the reconstruction error