

# 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

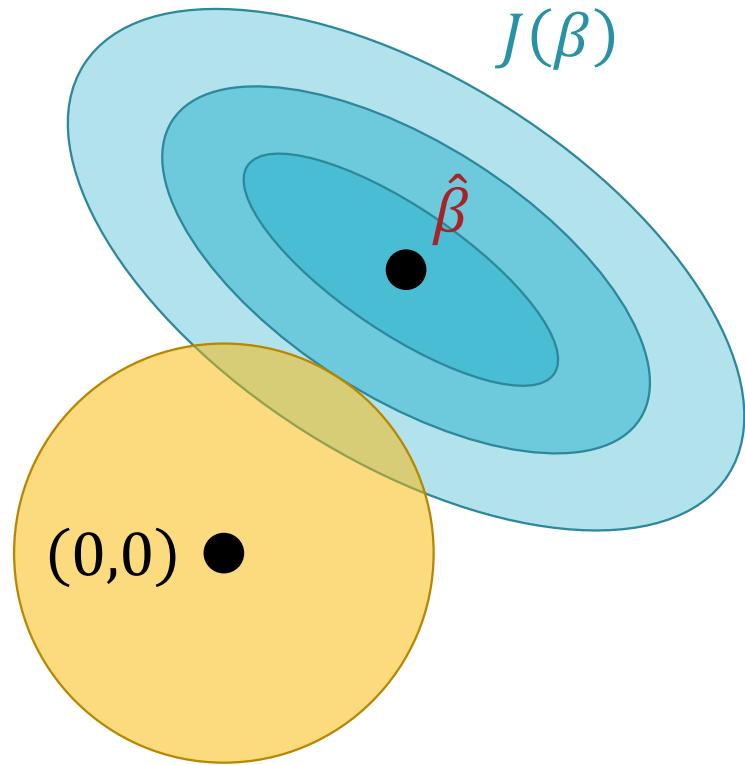
- Supervised learning -  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$ 
  - Regression -  $y^{(n)} \in \mathbb{R}$
  - Classification -  $y^{(n)} \in \{1, \dots, C\}$
- Unsupervised learning -  $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$ 
  - Clustering
  - Dimensionality reduction

# Learning Paradigms

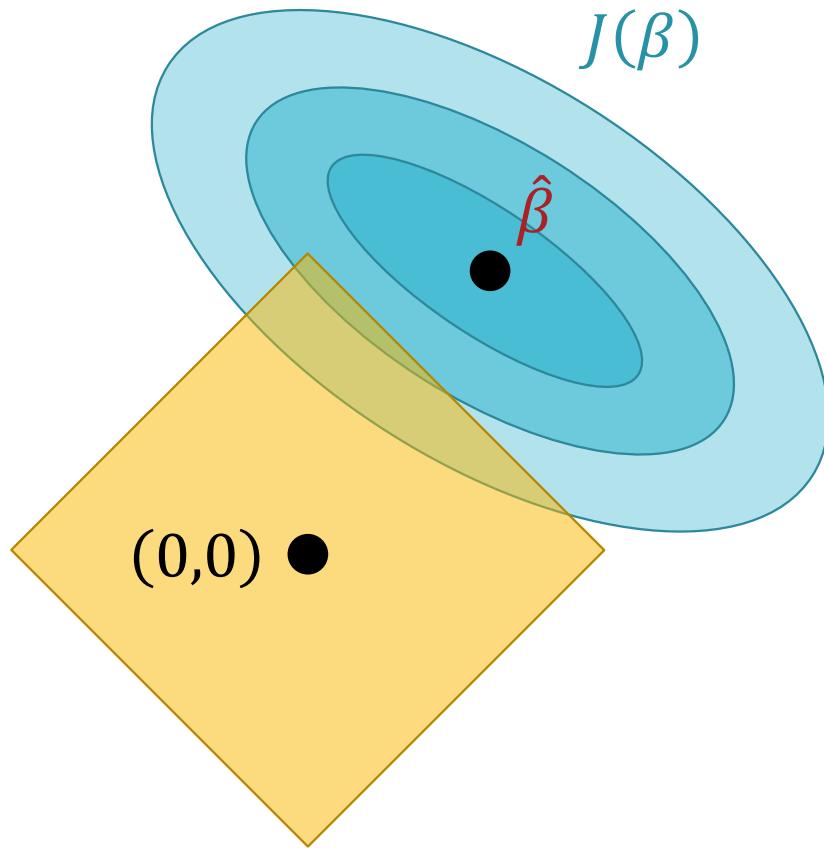
- Supervised learning -  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$ 
  - Regression -  $y^{(n)} \in \mathbb{R}$
  - Classification -  $y^{(n)} \in \{1, \dots, C\}$
- Unsupervised learning -  $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$ 
  - Clustering
  - Dimensionality reduction

# Dimensionality Reduction

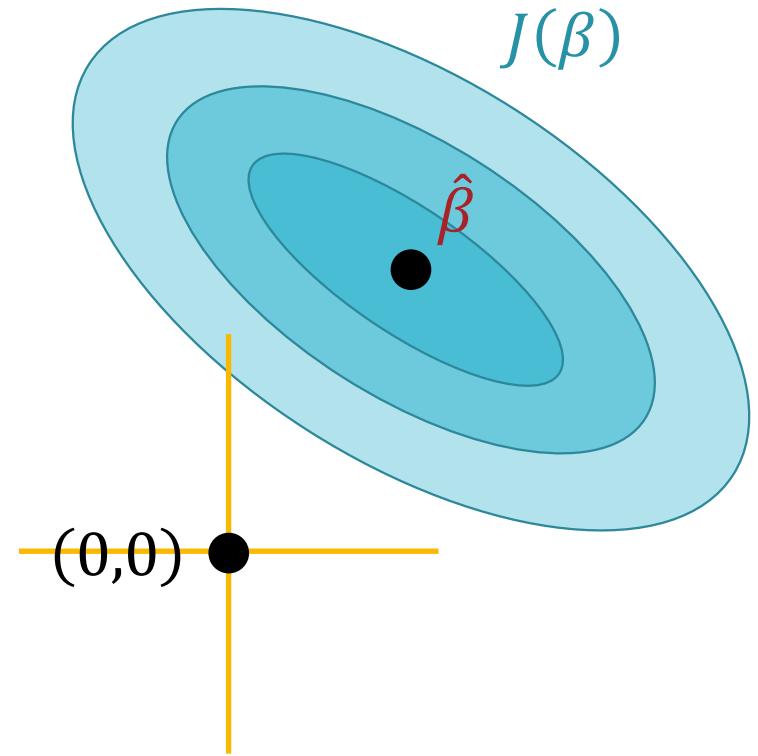
- Goal: given some unlabeled data set, learn a latent (typically lower-dimensional) representation
- Use cases:
  - Reducing computational cost (runtime, storage, etc...)
  - Improving generalization
  - Visualizing data
- Applications:
  - comparing images on the internet
  - demographic data
  - geospatial data



Ridge or  $L2$

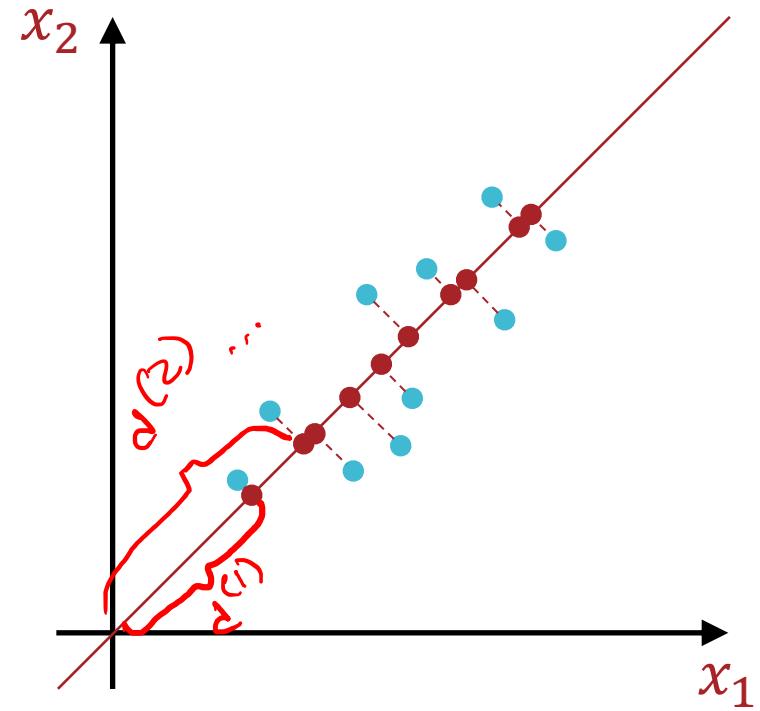
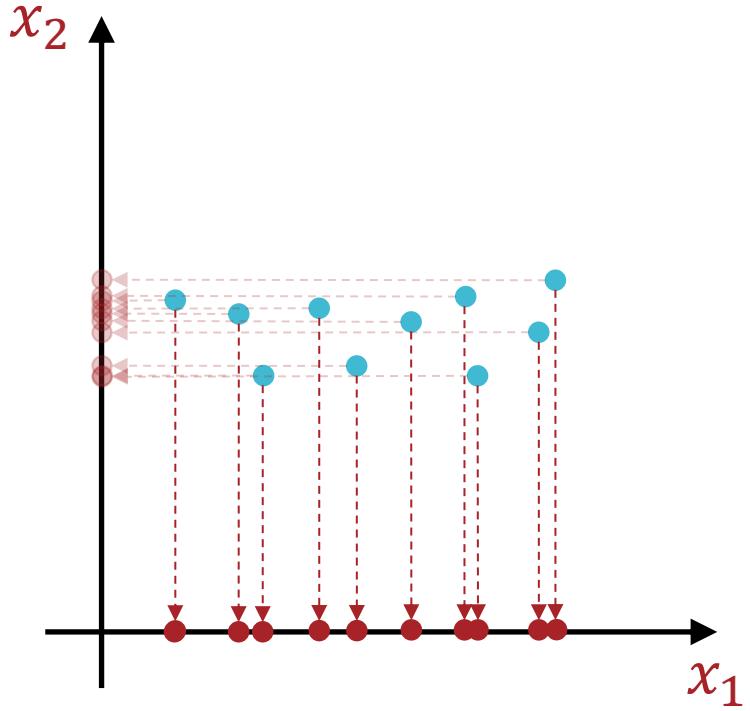


Lasso or  $L1$

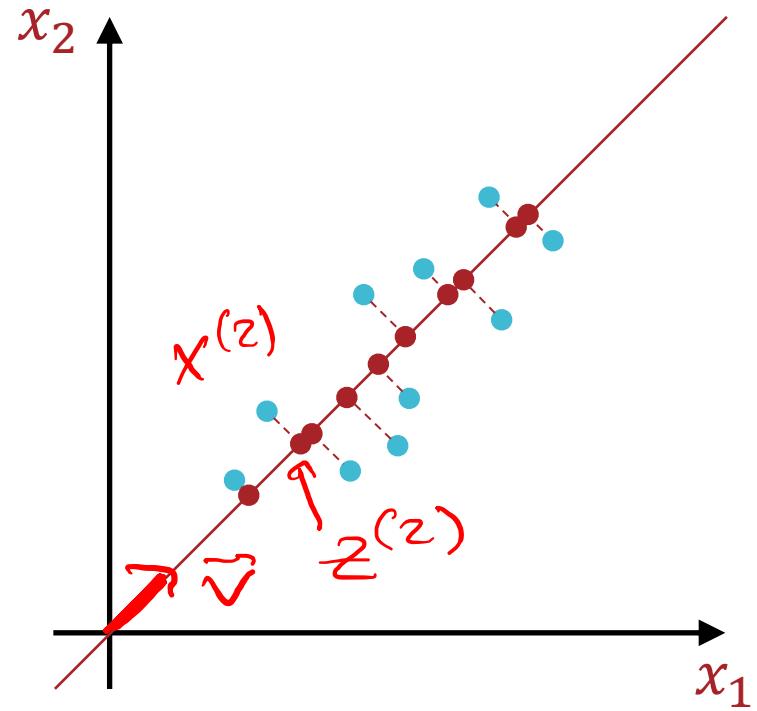
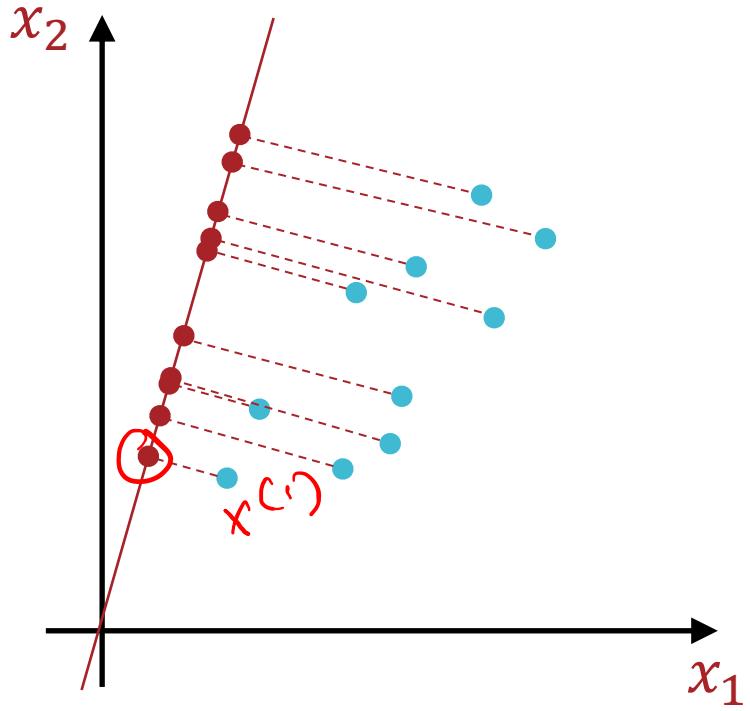


$L0$

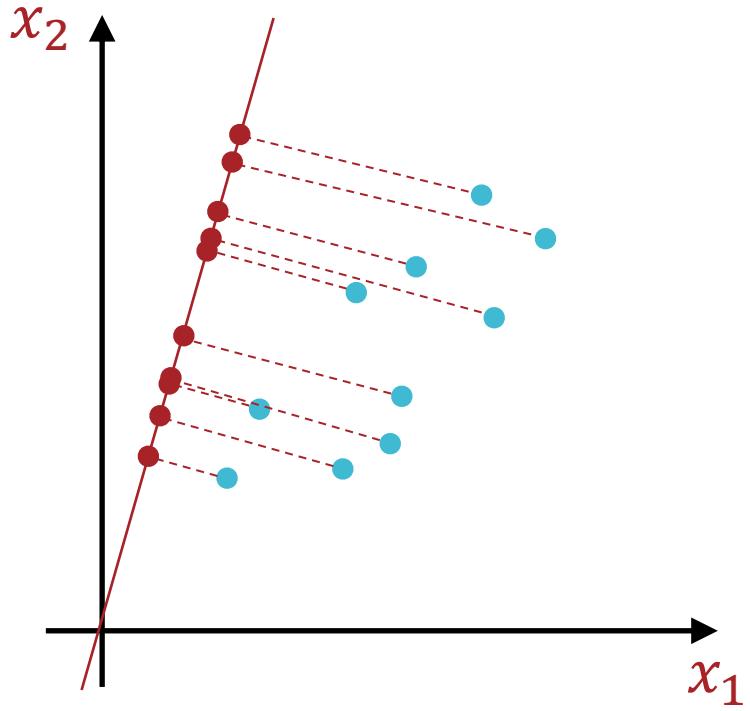
Recall:  $L1$  (or  $L0$ ) Regularization



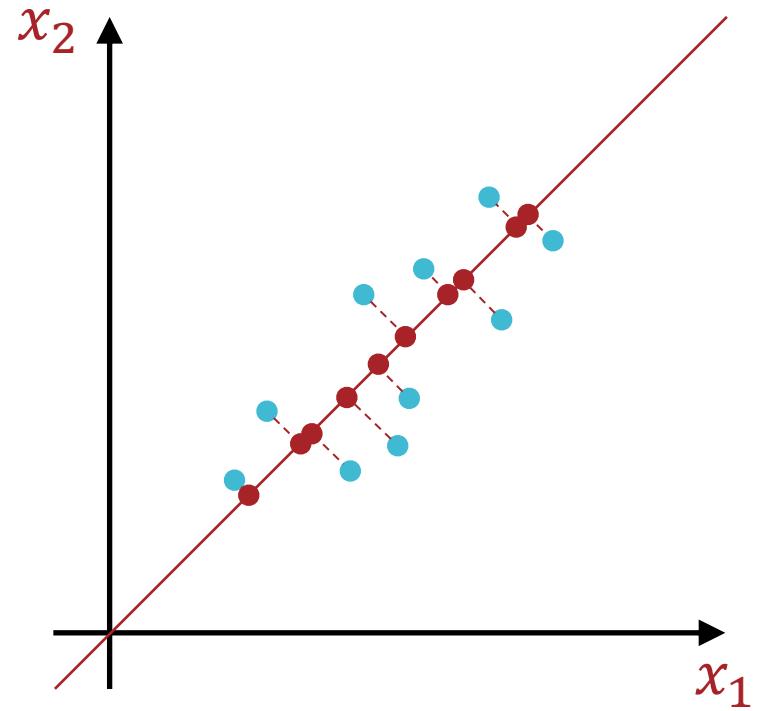
# Feature Elimination



# Feature Reduction



Option A



Option B

Which do you prefer *and why?*

**0 surveys completed**



**0 surveys underway**

# Which projection do you prefer?

---

Option A

Option B

Start the presentation to see live content. For screen share software, share the entire screen. Get help at [pollev.com/app](https://pollev.com/app)

# 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:

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

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

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

# Reconstruction Error

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

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

Length of projection

Direction of projection

# Reconstruction Error

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

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

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

$$\begin{aligned} & \left\| \tilde{x}^{(n)} - (v^T \tilde{x}^{(n)}) v \right\|_2^2 = (\tilde{x}^{(n)} - (v^T \tilde{x}^{(n)}) v)^T \\ & \quad (\tilde{x}^{(n)} - (v^T \tilde{x}^{(n)}) v) \\ & = \tilde{x}^{(n)T} \tilde{x}^{(n)} - 2(v^T \tilde{x}^{(n)}) v^T \tilde{x}^{(n)} + (v^T \tilde{x}^{(n)})^2 v^T v \\ & = \tilde{x}^{(n)T} \tilde{x}^{(n)} - (v^T \tilde{x}^{(n)})^2 \end{aligned}$$

# Minimizing the Reconstruction Error

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

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

$$\hat{v} = \underset{v: \|v\|_2^2=1}{\operatorname{argmax}} \sum_{n=1}^N (v^T \tilde{x}^{(n)})^2 = \sum_{n=1}^N (v^T \tilde{x}^{(n)}) (\tilde{x}^{(n)T} v)$$

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

## Maximizing the Variance

$$\|\vec{v}\|_2^2 = \vec{v}^\top \vec{v} = 1 \Rightarrow \vec{v}^\top \vec{v} - 1 = 0$$
$$\hat{\vec{v}} = \underset{\vec{v}: \|\vec{v}\|_2^2 = 1}{\operatorname{argmax}} \vec{v}^\top (X^\top X) \vec{v}$$
$$L(\vec{v}, \lambda) = \vec{v}^\top X^\top X \vec{v} - \lambda (\vec{v}^\top \vec{v} - 1)$$
$$\frac{\partial L}{\partial \vec{v}} = 2(X^\top X)\vec{v} - 2\lambda \vec{v}$$
$$\Rightarrow 2(X^\top X)\overset{\uparrow}{\vec{v}} - 2\lambda \overset{\uparrow}{\vec{v}} = 0$$
$$\Rightarrow (X^\top X)\overset{\uparrow}{\vec{v}} = \lambda \overset{\uparrow}{\vec{v}}$$
$$\Rightarrow \overset{\uparrow}{\vec{v}} \text{ is an eigenvector of } X^\top X$$

but which one?

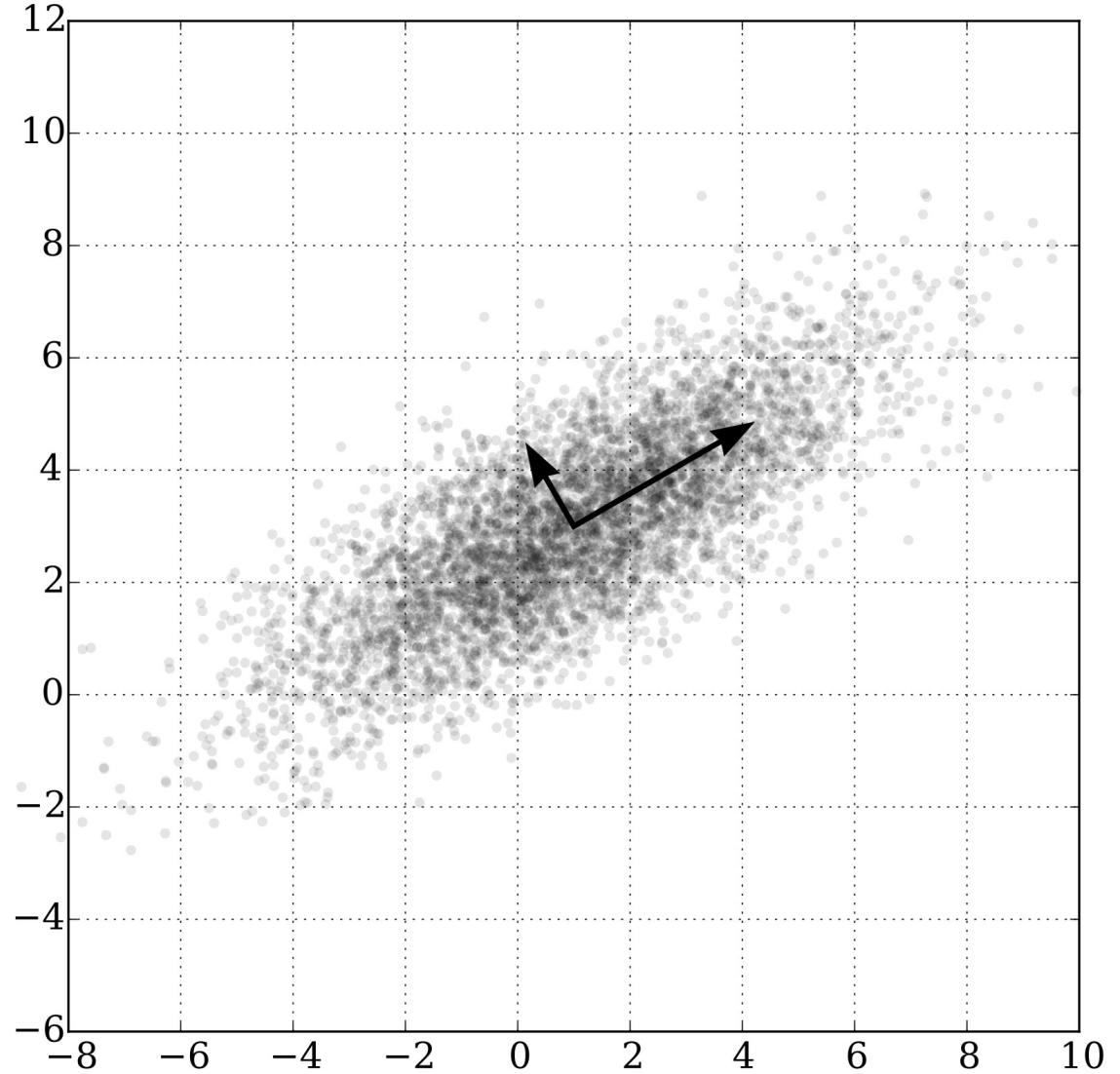
# Maximizing the Variance

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

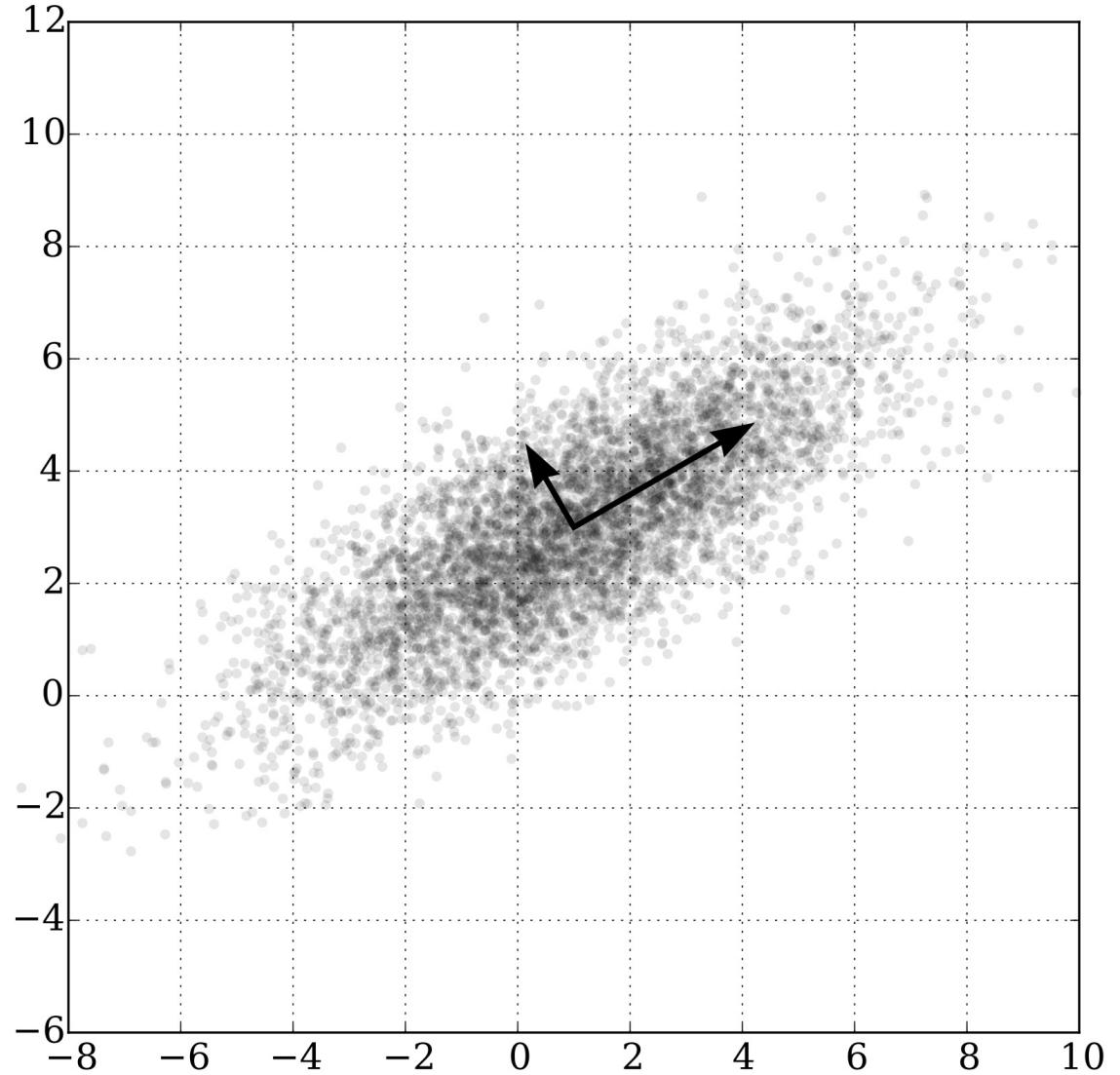
$$(X^T X) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}} \rightarrow \hat{\mathbf{v}}^T (X^T X) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}}^T \hat{\mathbf{v}} = \lambda$$

- 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  $X^TX$
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  $X^TX$

# 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  $X^T X$
  3. Collect the top  $\rho$  eigenvectors (corresponding to the  $\rho$  largest eigenvalues),  $V_\rho \in \mathbb{R}^{D \times \rho}$
  4. Project the data into the space defined by  $V_\rho$ ,  $Z = X V_\rho$
- Output:  $Z$ , the transformed (potentially lower-dimensional) data

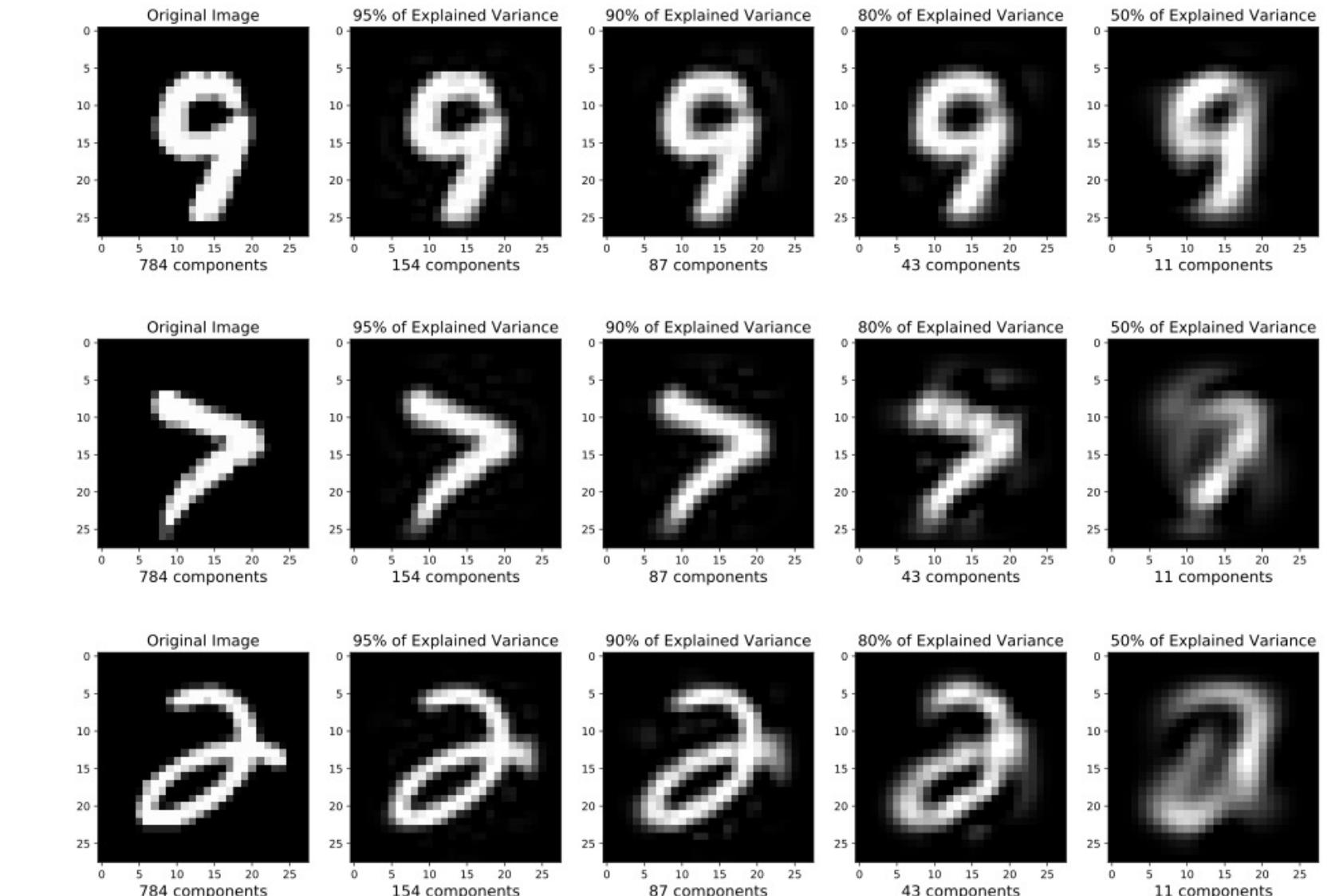
# How many PCs should we use?

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

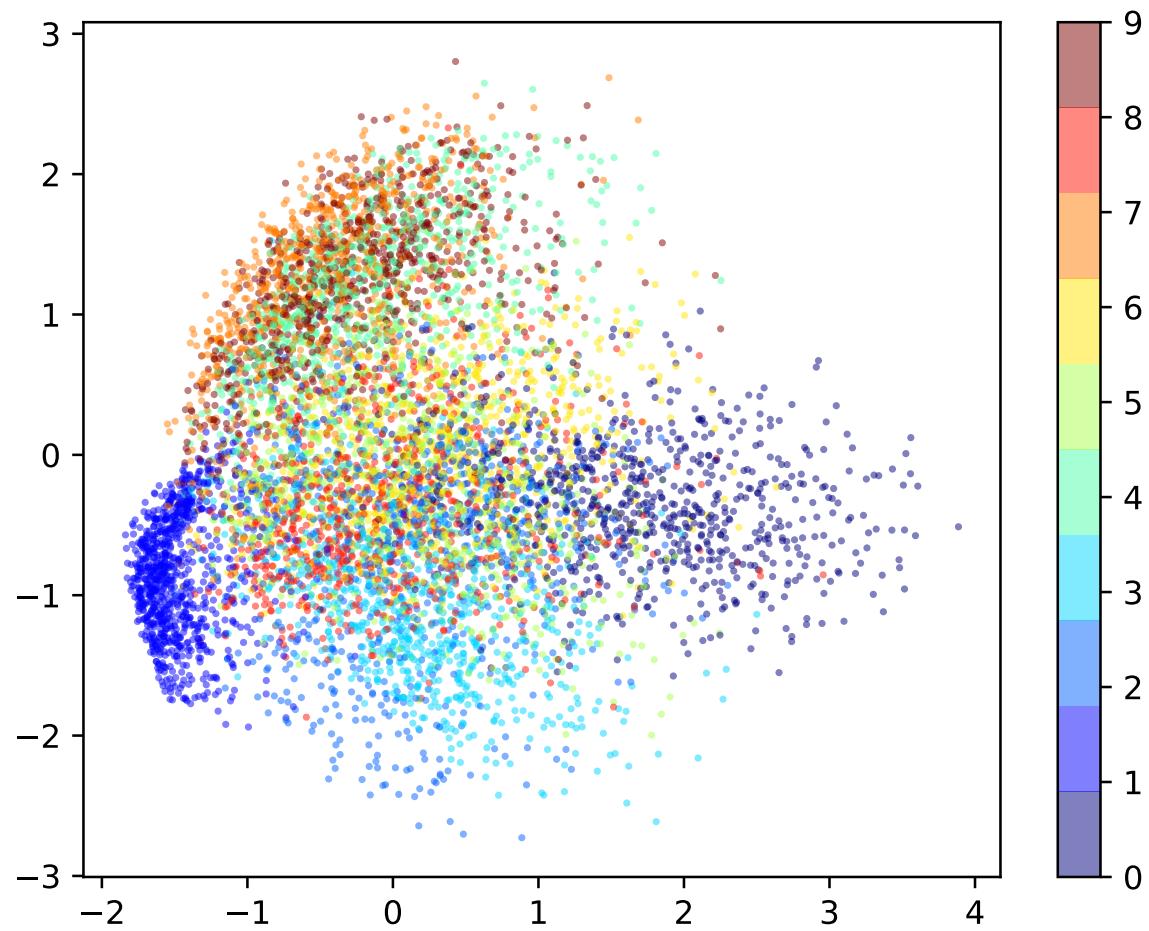
# Choosing the number of PCs

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

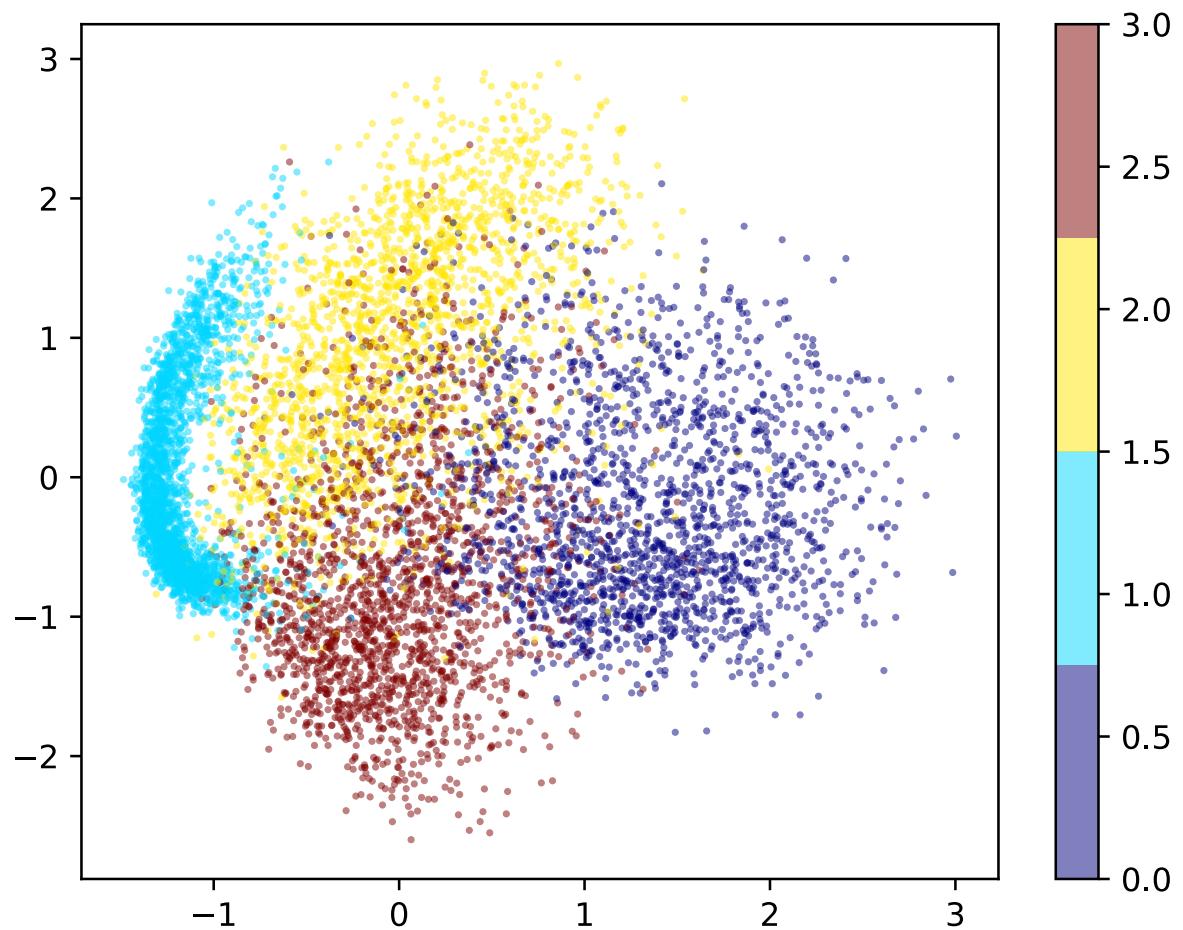
# PCA Example: MNIST Digits



# 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