

10-301/601: Introduction to Machine Learning

Lecture 18: Dimensionality Reduction

Henry Chai

7/10/23

Front Matter

- Announcements
 - PA4 released 6/15, due 7/13 at 11:59 PM
 - Quiz 6: Deep Learning & Learning Theory on 7/11 (tomorrow!)
- Recommended Readings
 - Murphy, [Chapters 12.2.1 - 12.2.3](#)
 - Daumé III, [Chapter 15: Unsupervised Learning](#)

Learning Paradigms

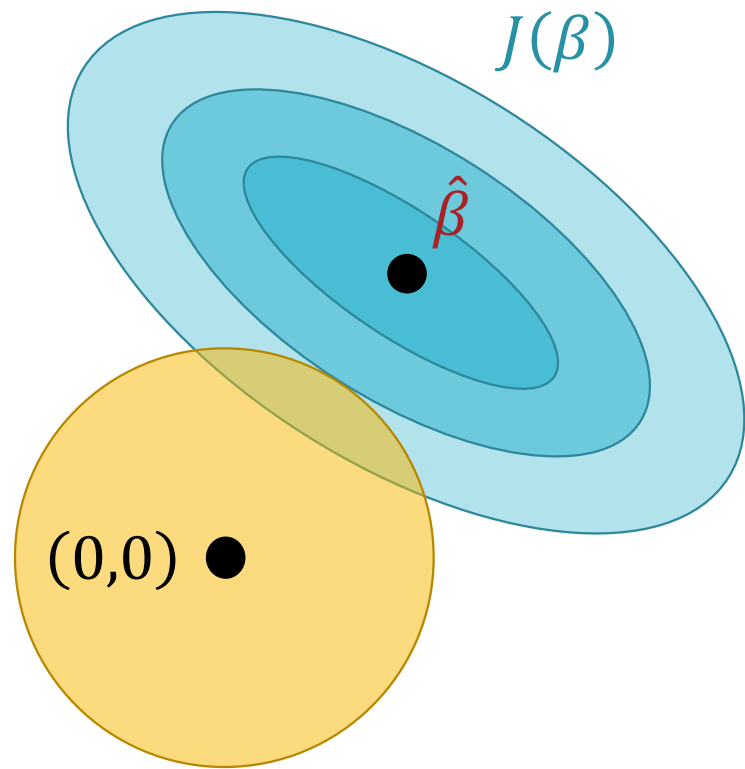
- Supervised learning - $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$
 - Regression - $y^{(n)} \in \mathbb{R}$
 - Classification - $y^{(n)} \in \{1, \dots, C\}$
- Unsupervised learning - $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$
 - Clustering
 - Dimensionality reduction

Learning Paradigms

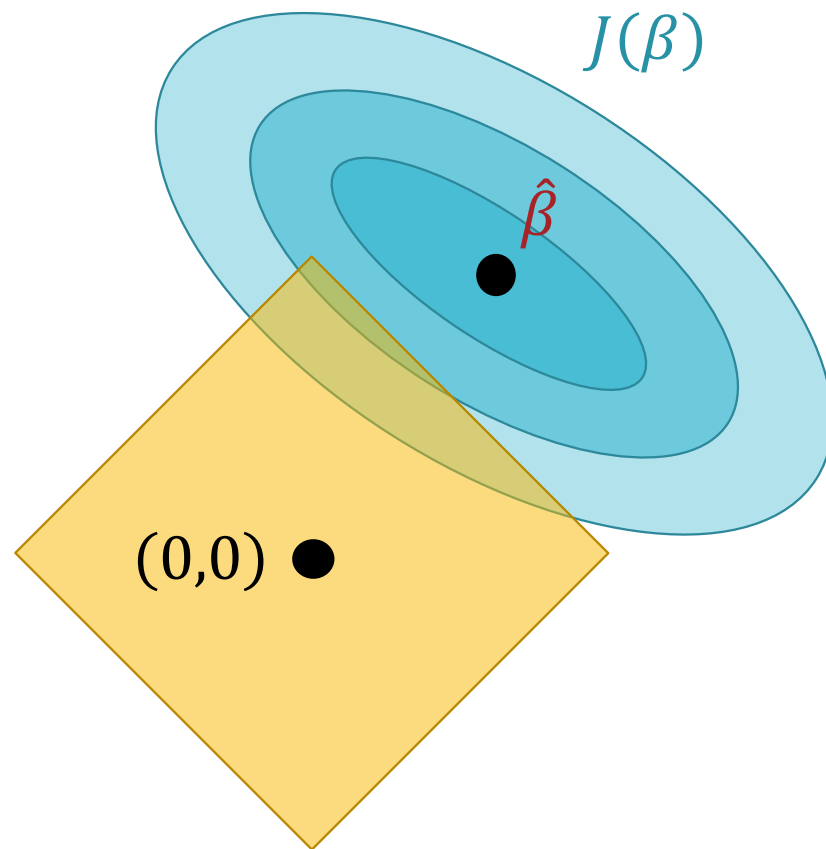
- Supervised learning - $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$
 - Regression - $y^{(n)} \in \mathbb{R}$
 - Classification - $y^{(n)} \in \{1, \dots, C\}$
- Unsupervised learning - $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$
 - Clustering
 - **Dimensionality reduction**

Dimensionality Reduction

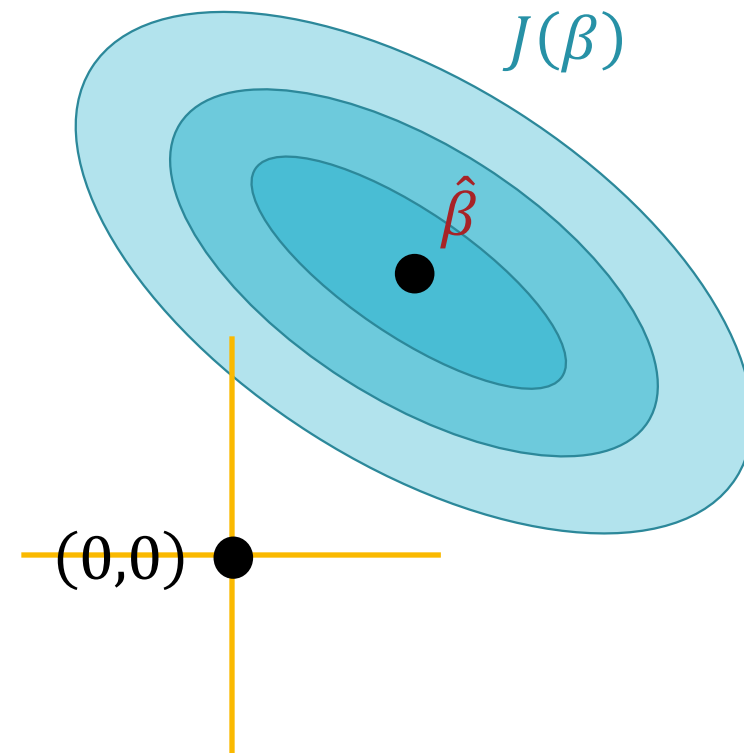
- Goal: given some unlabeled data set, learn a latent (typically lower-dimensional) representation
- Use cases:
 - Reducing computational cost (runtime, storage, etc...)
 - Improving generalization
 - Visualizing data
- Applications:
 - Recommender system: demographic info., usage behavior, current recommendations
 - Inspection of parts: physical measurements of components - most are typically irrelevant
 - Text & image/video



Ridge or L_2

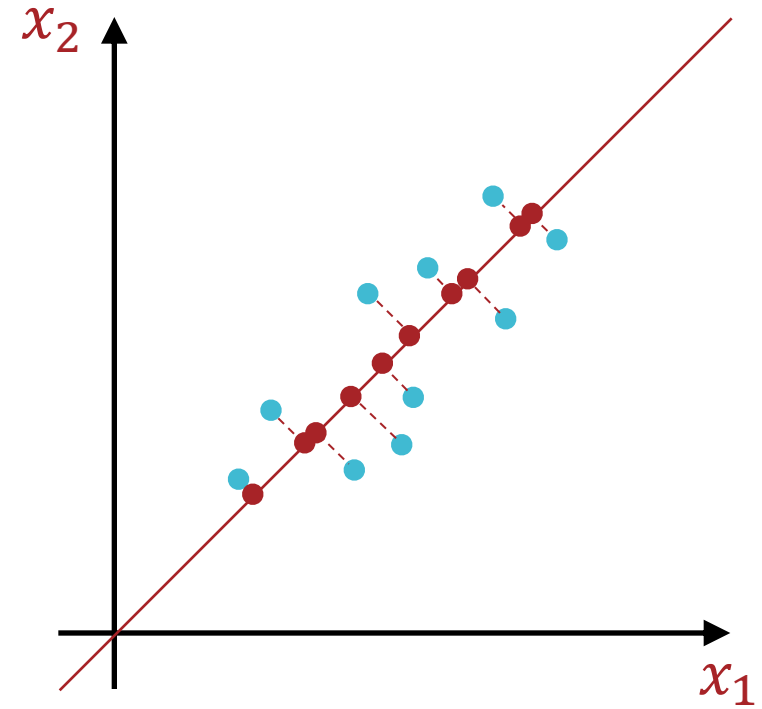
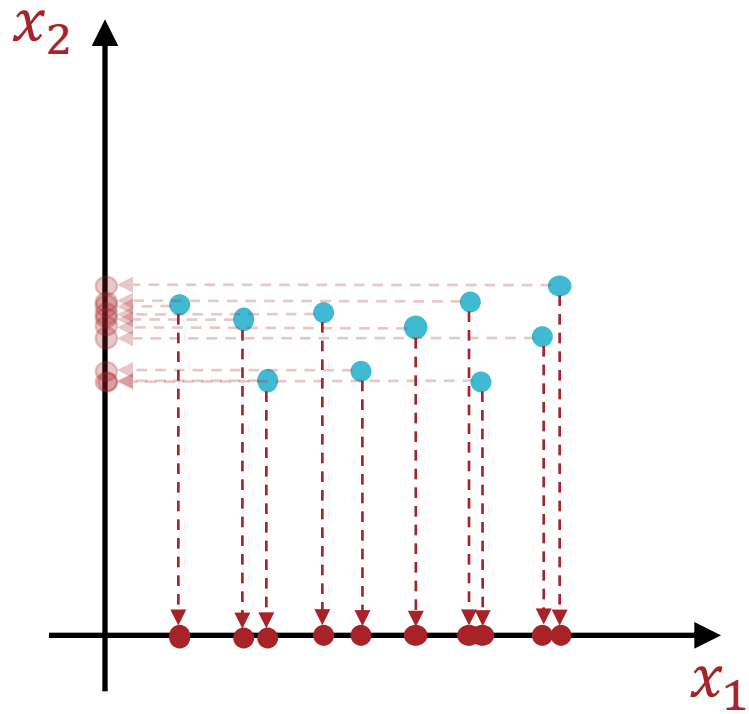


Lasso or L_1

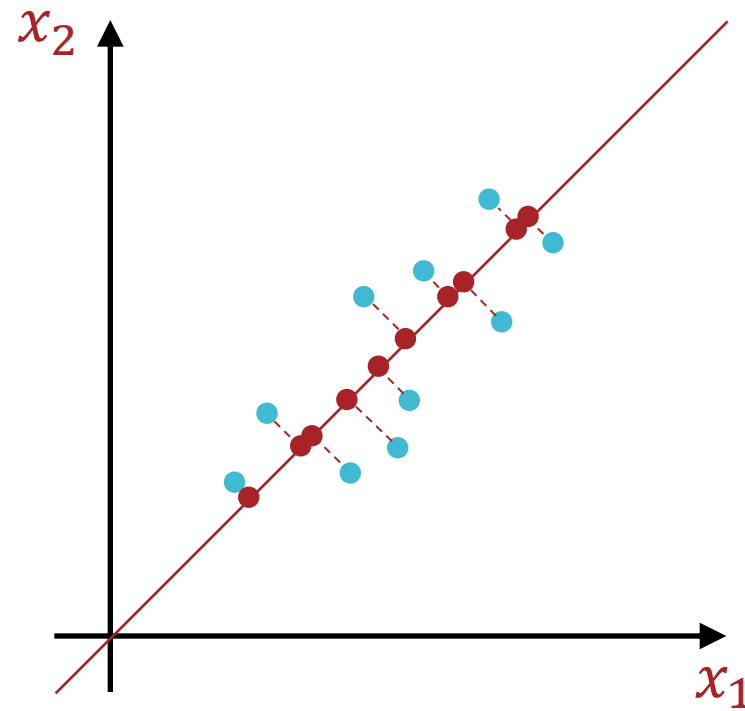
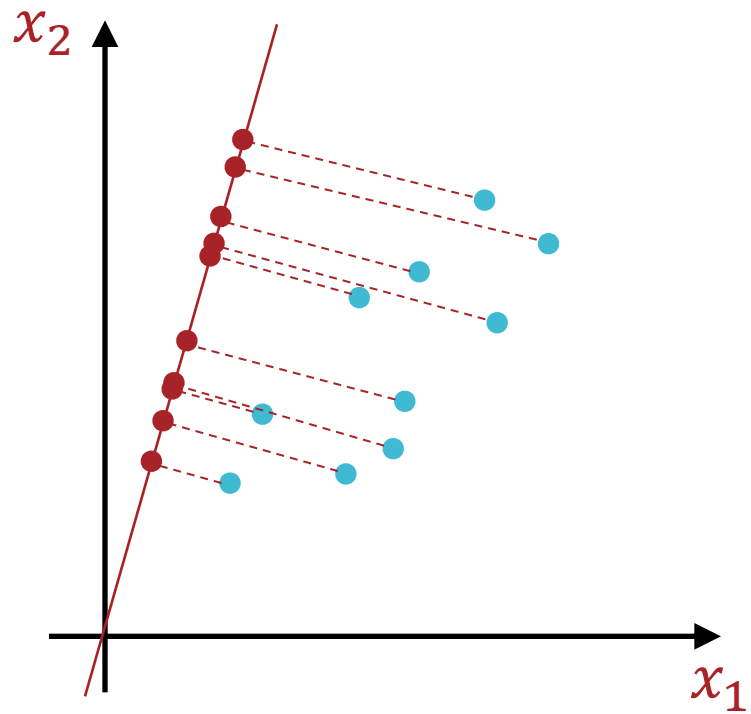


L_0

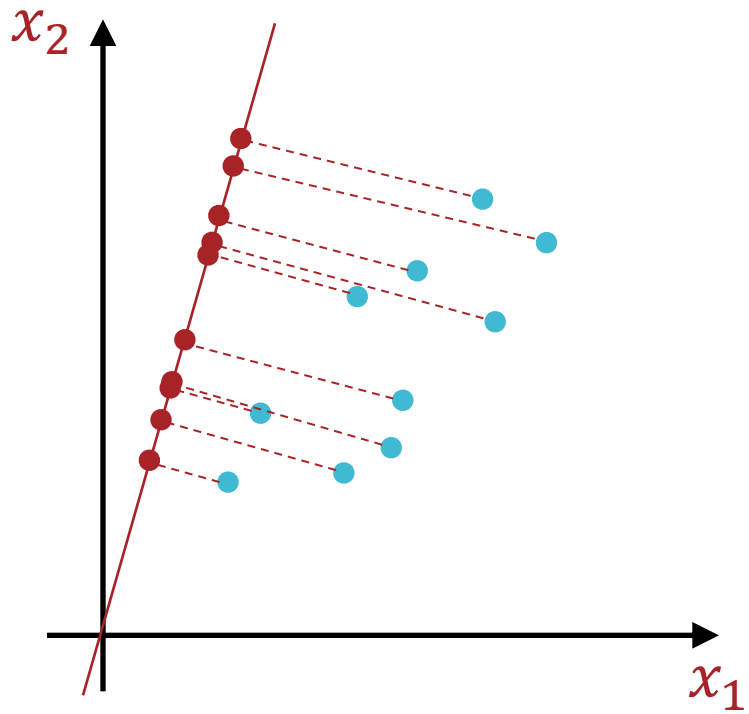
Recall: L_1 (or L_0) Regularization



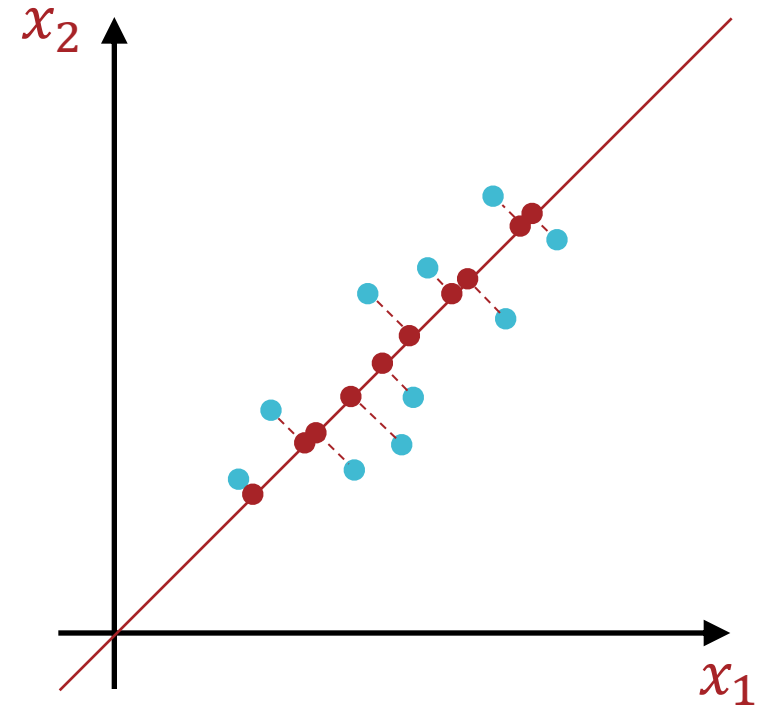
Feature Elimination



Feature Reduction



Option A



Option B

Which projection do you prefer?

Which projection do you prefer?

Option
A

Option
B

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:

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

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

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

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

$$v^T \tilde{x}^{(n)} \in \mathbb{R}$$

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

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

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

$$\begin{aligned} \|\tilde{x}^{(n)} - (v^T \tilde{x}^{(n)})v\|_2^2 &= (\tilde{x}^{(n)} - (v^T \tilde{x}^{(n)})v)^T (\tilde{x}^{(n)} - (v^T \tilde{x}^{(n)})v) \\ &= \tilde{x}^{(n)T} \tilde{x}^{(n)} - 2(v^T \tilde{x}^{(n)})v^T \tilde{x}^{(n)} + ((v^T \tilde{x}^{(n)})v)^T ((v^T \tilde{x}^{(n)})v) \\ &= \tilde{x}^{(n)T} \tilde{x}^{(n)} - 2(v^T \tilde{x}^{(n)})(v^T \tilde{x}^{(n)}) + (v^T \tilde{x}^{(n)})^2 \underbrace{v^T v}_1 \\ &= \tilde{x}^{(n)T} \tilde{x}^{(n)} - 2(v^T \tilde{x}^{(n)})^2 + (v^T \tilde{x}^{(n)})^2 \\ &= \|\tilde{x}^{(n)}\|_2^2 - (v^T \tilde{x}^{(n)})^2 \end{aligned}$$

Reconstruction Error

Minimizing the Reconstruction Error



Maximizing the Variance

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

$$\hat{v} = \underset{v: \|v\|_2=1}{\operatorname{argmin}} \sum_{n=1}^N \left(\underbrace{\|\tilde{x}^{(n)}\|_2^2}_{\text{variance of a zero-mean random variable}} - (v^T \tilde{x}^{(n)})^2 \right)$$

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

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

$$\hat{v} = \underset{v: \|v\|_2=1}{\operatorname{argmax}} \frac{1}{N} \sum_{n=1}^N (v^T \tilde{x}^{(n)})^2 v^T v - 0$$

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

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

Maximizing the Variance

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

$$\begin{aligned} \rightarrow \text{s.t. } \|v\|_2^2 &= 1 \\ &\rightarrow \|v\|_2^2 - 1 = 0 \end{aligned}$$

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

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

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

$\rightarrow \hat{v}$ is an ~~eigenvector~~ eigenvector of $X^T X$!

but which one?

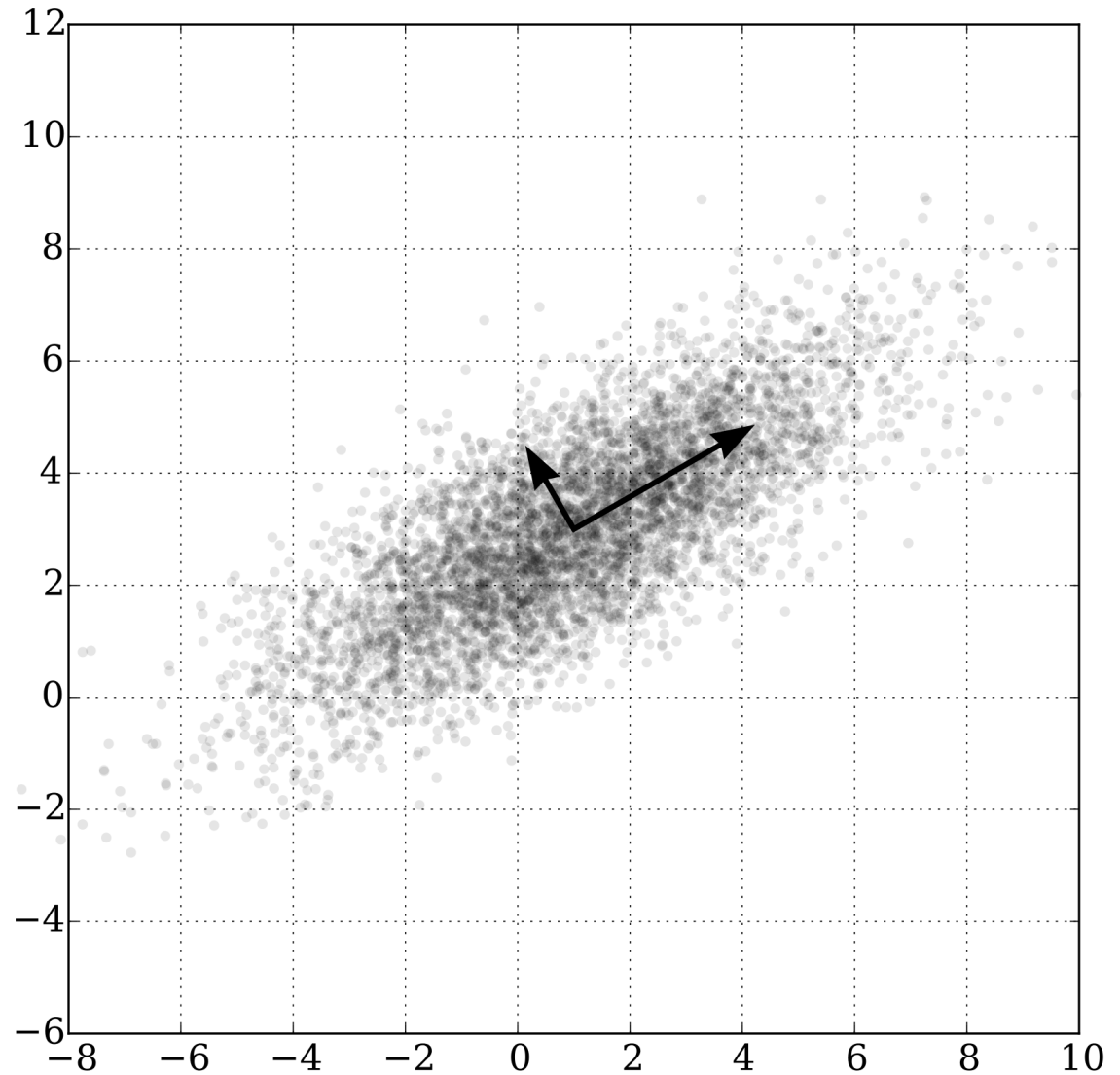
Maximizing the Variance

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

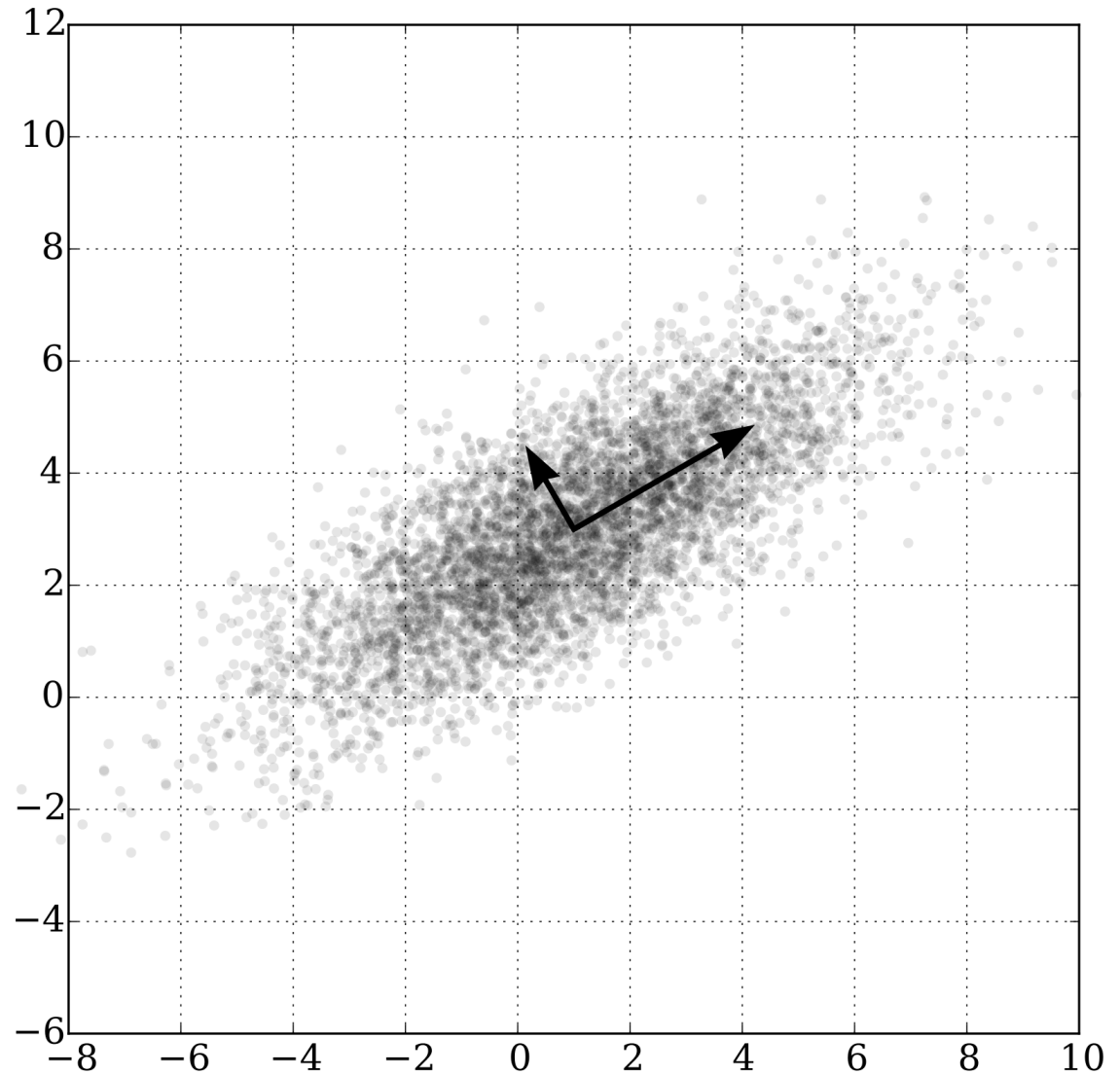
$$(\mathbf{X}^T \mathbf{X}) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}} \rightarrow \hat{\mathbf{v}}^T (\mathbf{X}^T \mathbf{X}) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}}^T \hat{\mathbf{v}} = \lambda$$

- The first principal component is the eigenvector $\hat{\mathbf{v}}_1$ that corresponds to the largest eigenvalue λ_1
- The second principal component is the eigenvector $\hat{\mathbf{v}}_2$ that corresponds to the second largest eigenvalue λ_2
 - $\hat{\mathbf{v}}_1$ and $\hat{\mathbf{v}}_2$ are orthogonal
- Etc ...
- λ_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 $X^T X$
 3. Collect the top ρ eigenvectors (corresponding to the ρ largest eigenvalues), $V_\rho \in \mathbb{R}^{D \times \rho}$
 4. Project the data into the space defined by V_ρ , $Z = X V_\rho$
- 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 ρ eigenvectors (corresponding to the ρ largest eigenvalues), $V_\rho \in \mathbb{R}^{D \times \rho}$
 4. Project the data into the space defined by V_ρ , $Z = X V_\rho$
- Output: Z , the transformed (potentially lower-dimensional) data

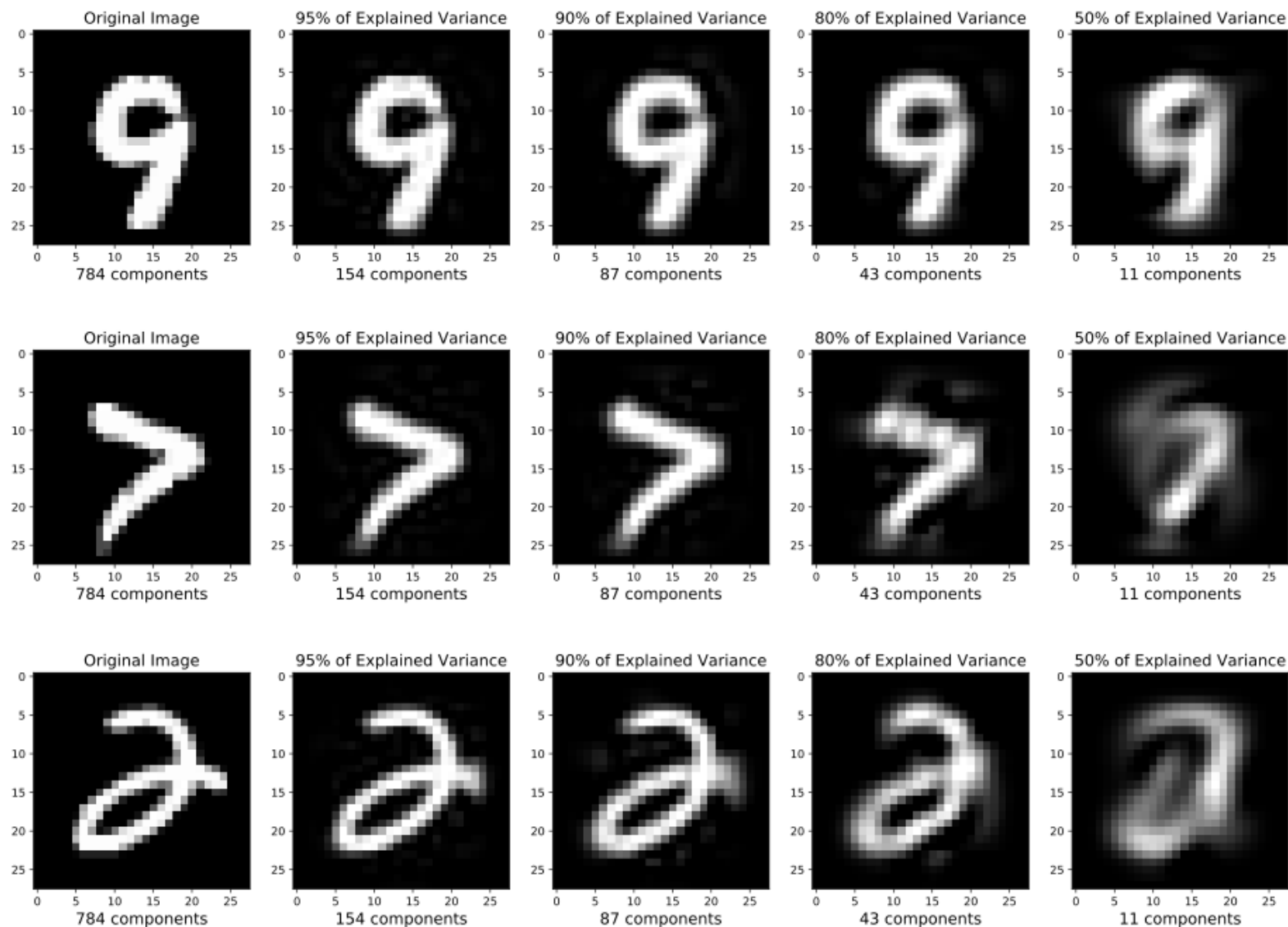
Choosing the number of PCs

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

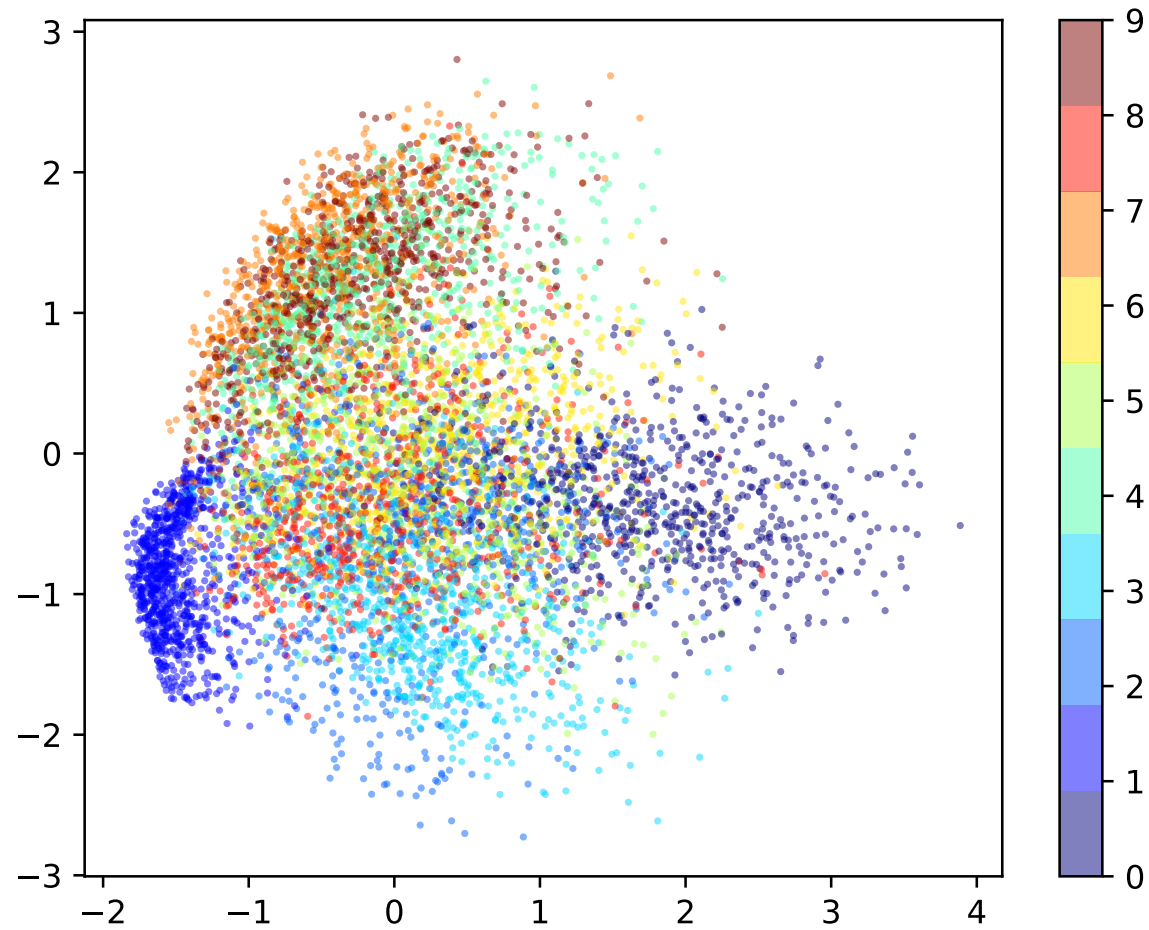
$$\lambda_i / \sum \lambda_j$$

- 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

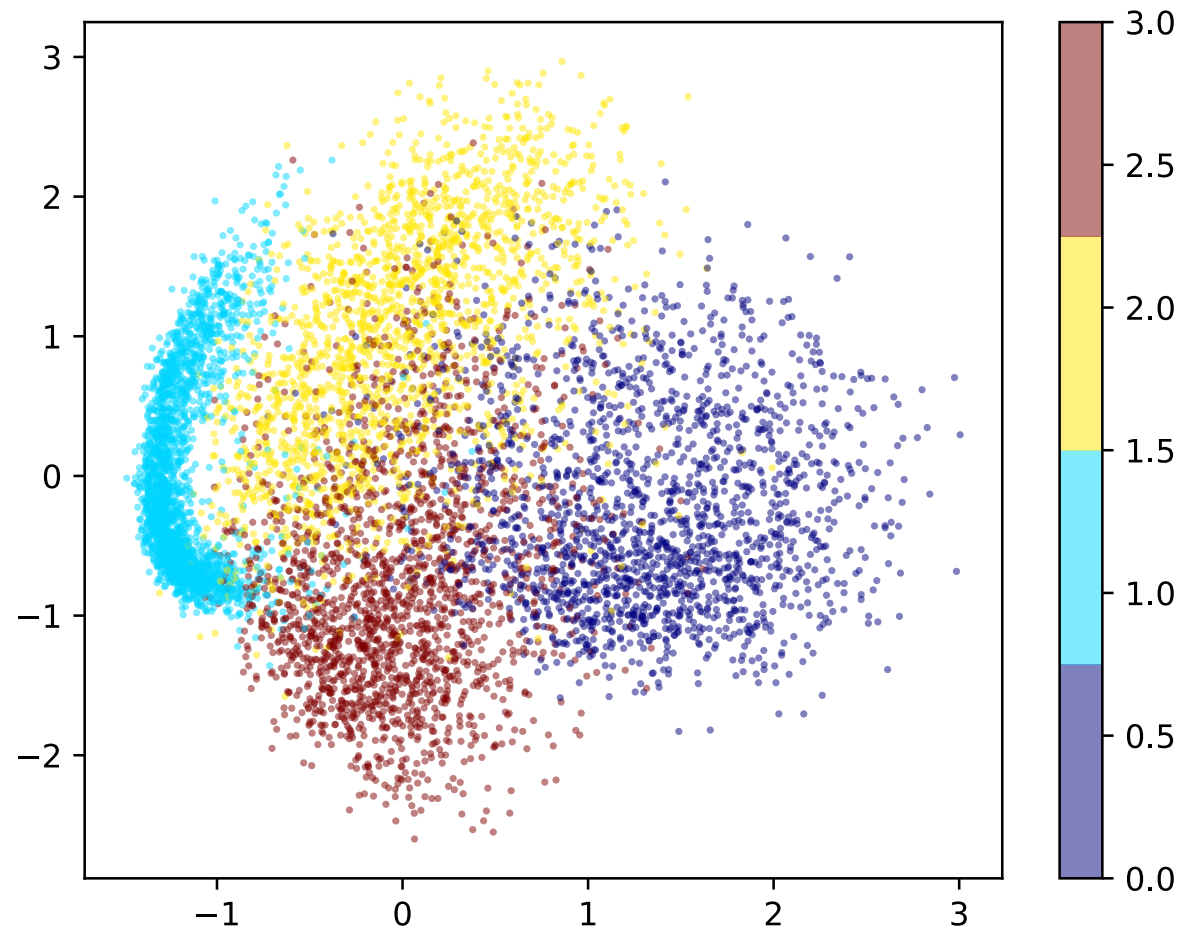
PCA Example: MNIST Digits



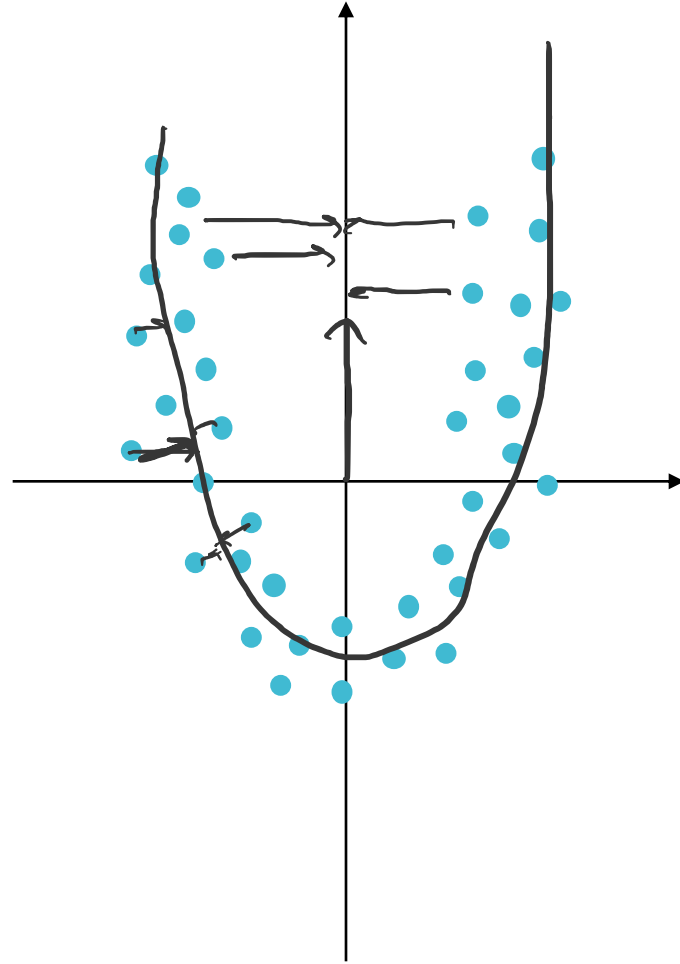
PCA Example: MNIST Digits



PCA Example: MNIST Digits

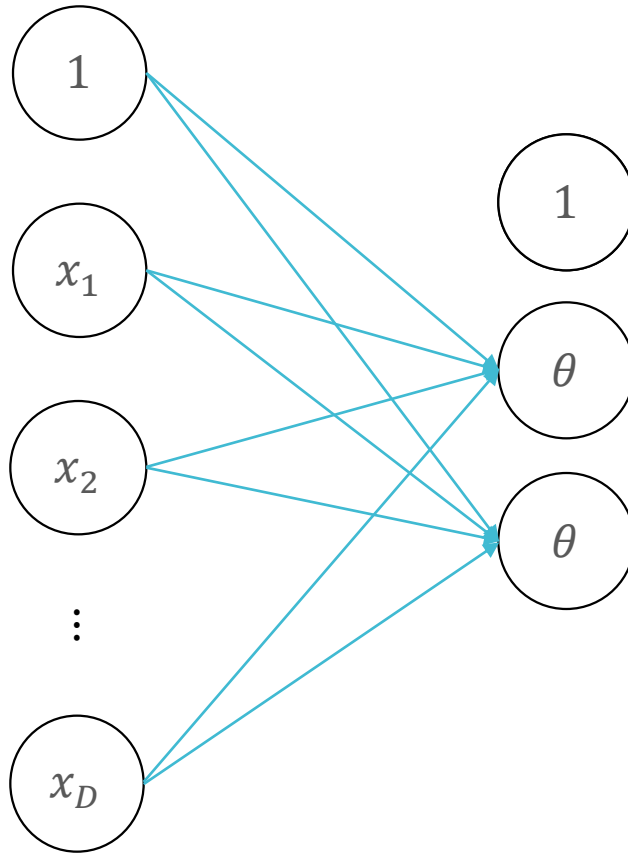


Shortcomings of PCA



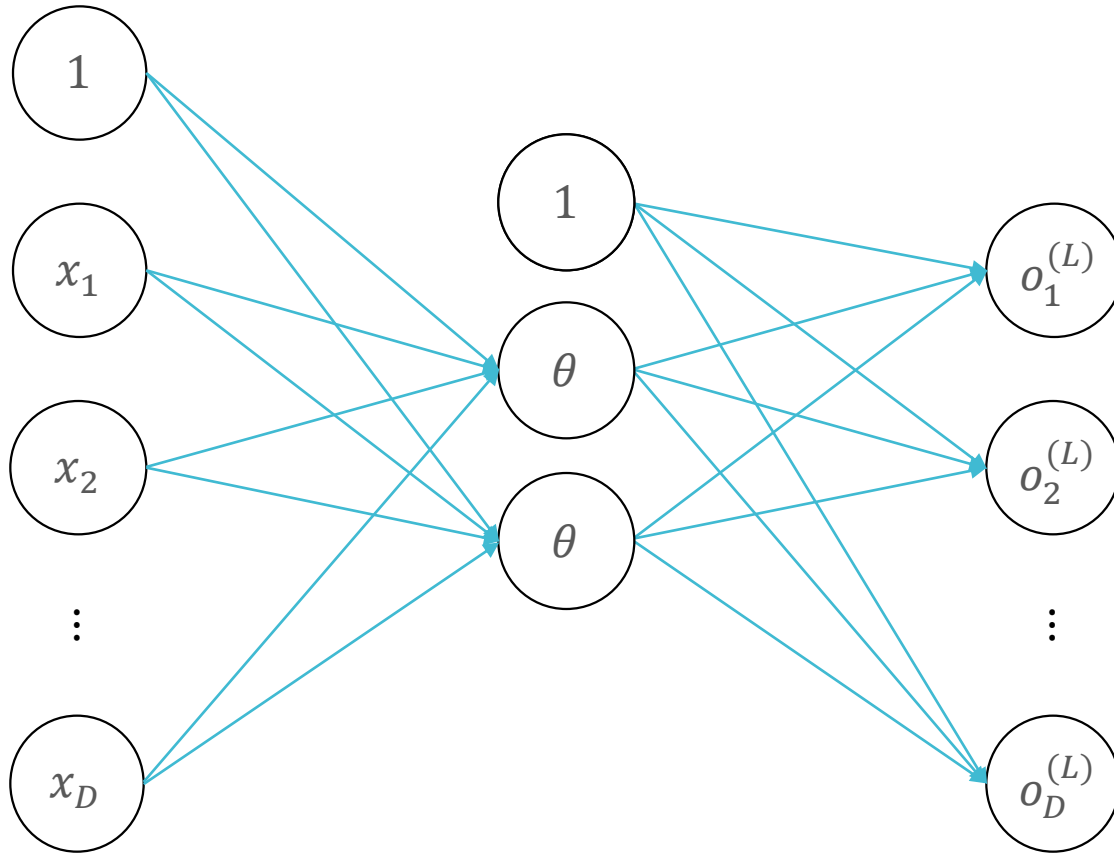
- PCs are
constrained to be
linear

Autoencoders



Insight: neural networks implicitly learn low-dimensional representations of inputs in hidden layers

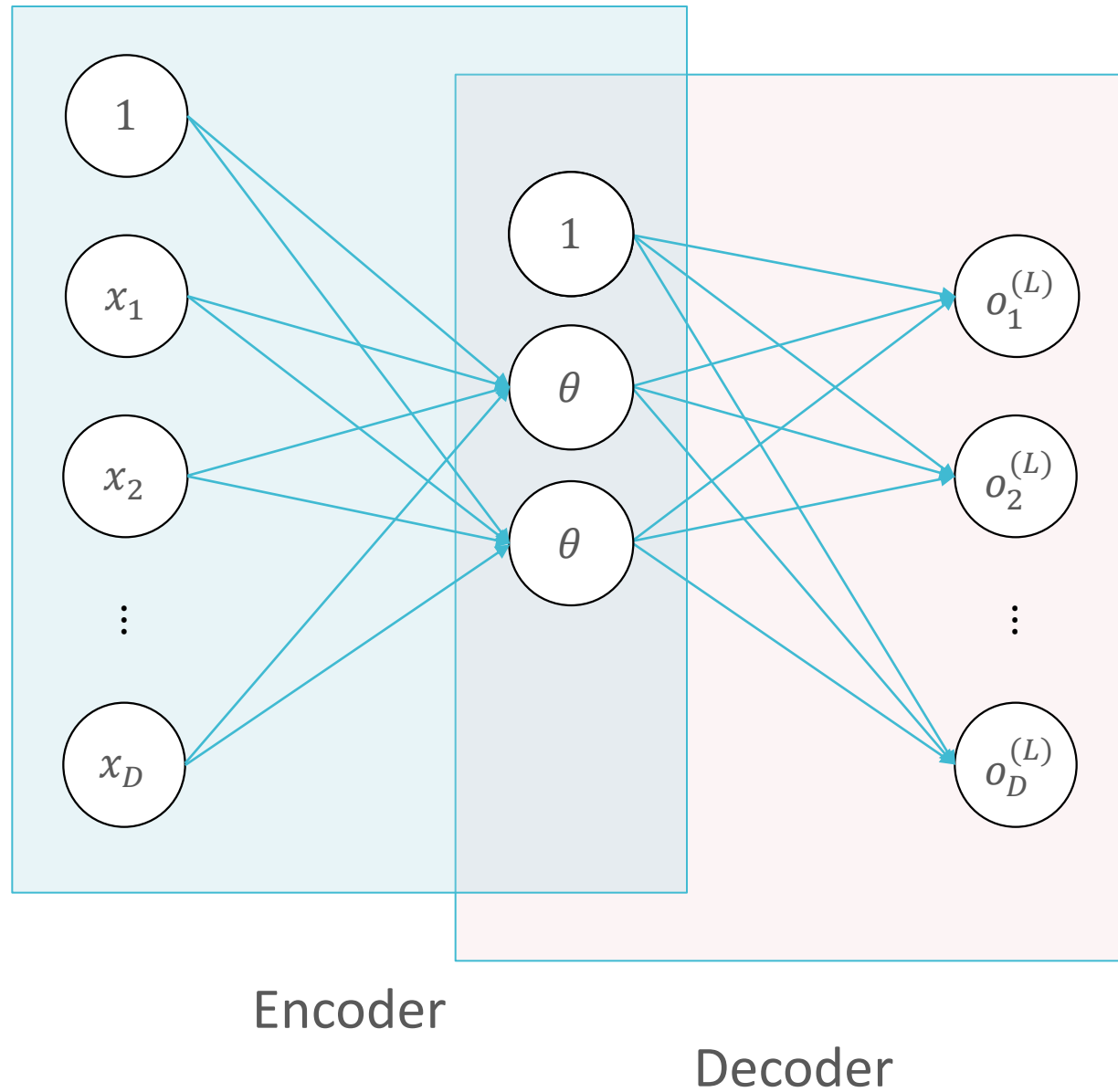
Autoencoders



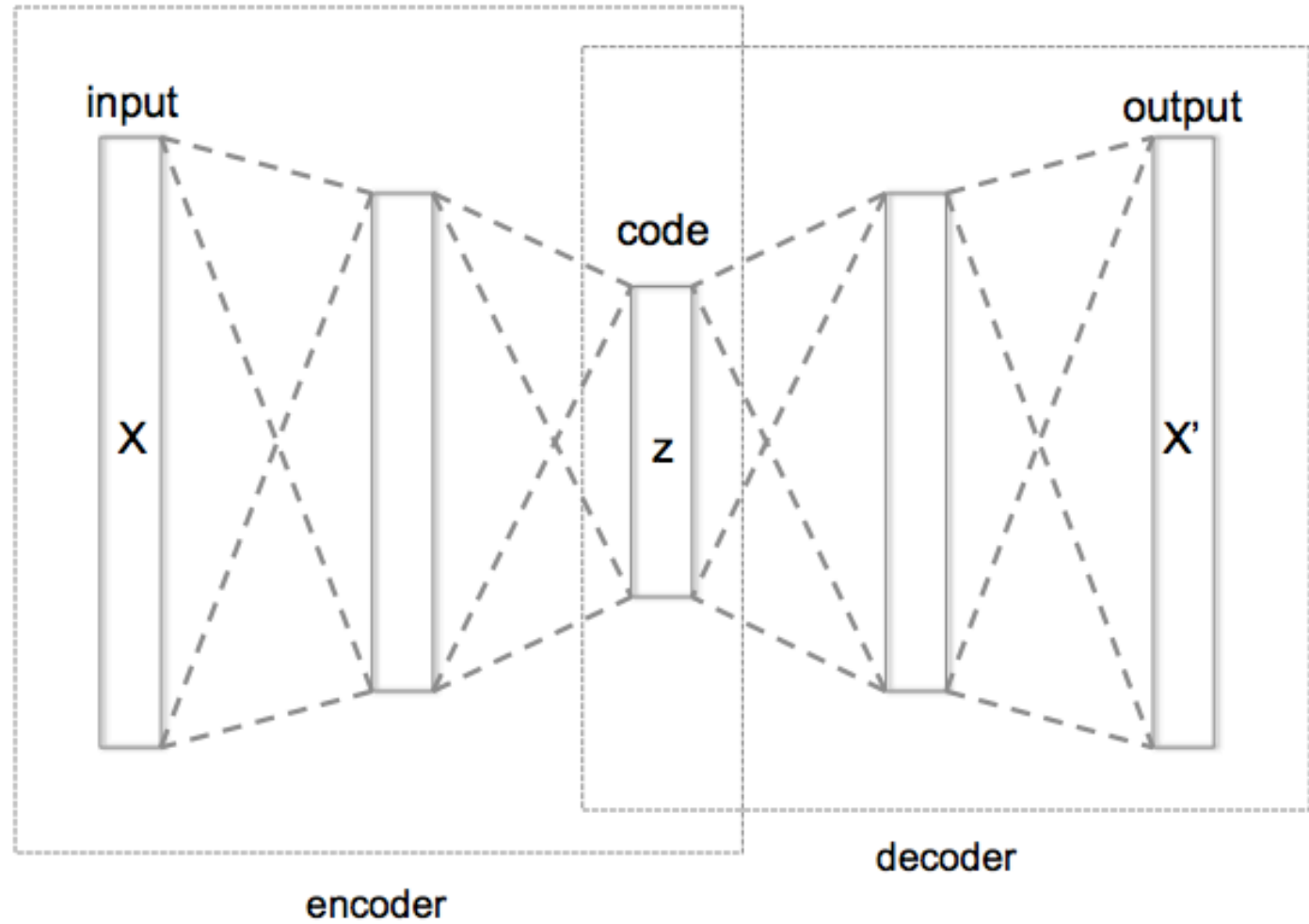
- Learn the weights by minimizing the reconstruction loss:

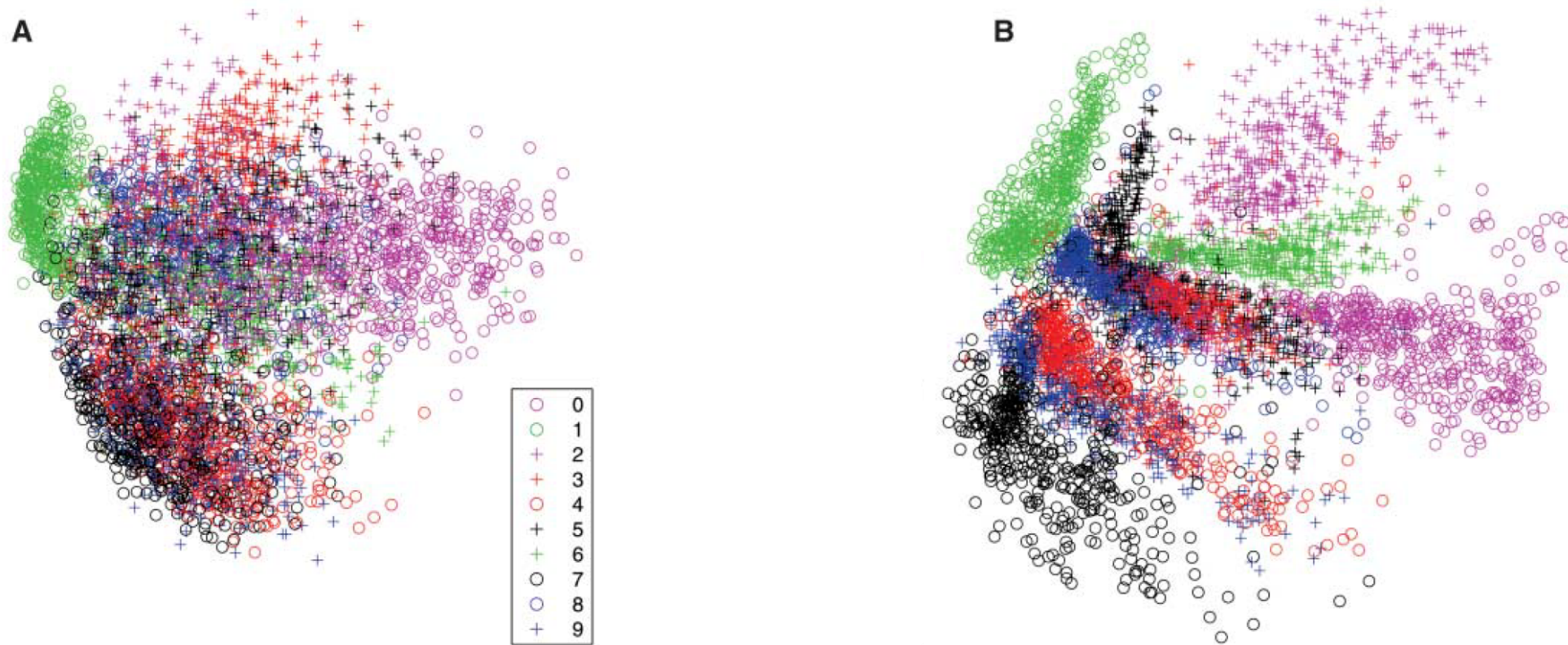
$$e(\mathbf{x}) = \|\mathbf{x} - \mathbf{o}^{(L)}\|_2^2$$

Autoencoders



Deep Autoencoders





PCA (A) vs. Autoencoders (B) (Hinton and Salakhutdinov, 2006)

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