

10-301/601: Introduction to Machine Learning

Lecture 25: Dimensionality Reduction

Matt Gormley & Henry Chai

11/25/24

Front Matter

- Announcements
 - HW8 released 11/17, due 11/25 (today!) at 11:59 PM
 - Please be mindful of your grace day usage (see [the course syllabus](#) for the policy)
 - HW9 released 11/25 (today!), due 12/5 at 11:59 PM
 - You are not expected to work on HW9 over Thanksgiving break
 - Relatedly, there are no OH over break
 - **You may only take at most 2 grace days on HW9**
 - Recitation for HW9 on Monday, 12/2 (after break)

Recall: K-means Algorithm

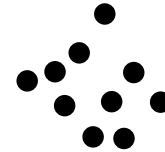
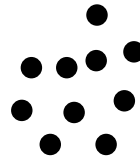
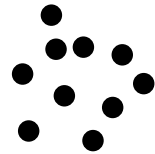
- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)})\}_{n=1}^N, K$
 1. Initialize cluster centers $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$
 2. While NOT CONVERGED
 - a. Assign each data point to the cluster with the nearest cluster center:
$$z^{(n)} = \underset{k}{\operatorname{argmin}} \|\mathbf{x}^{(n)} - \boldsymbol{\mu}_k\|_2$$
 - b. Recompute the cluster centers:
$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n: z^{(n)}=k} \mathbf{x}^{(n)}$$
where N_k is the number of data points in cluster k
- Output: cluster centers $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$ and cluster assignments $z^{(1)}, \dots, z^{(N)}$

Setting K

- Idea: choose the value of K that minimizes the objective function

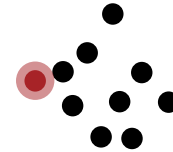
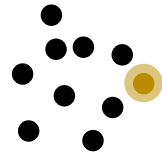
Initializing K -means: Lloyd's Method

- Randomly choose K data points to be the initial cluster centers



Initializing K -means: Lloyd's Method

- Randomly choose K data points to be the initial cluster centers



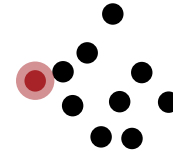
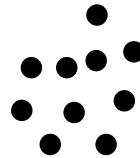
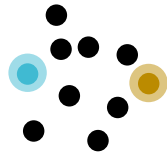
Initializing K -means: Lloyd's Method

- Randomly choose K data points to be the initial cluster centers



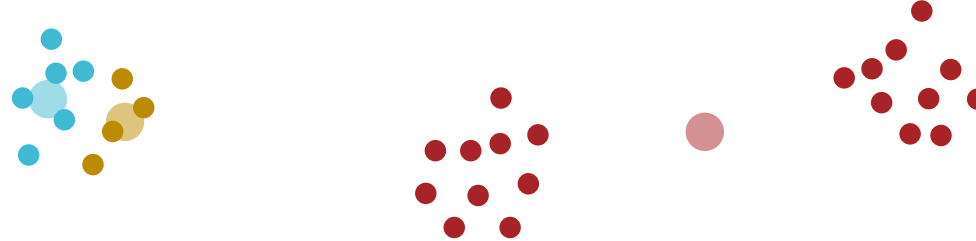
Initializing K -means: Lloyd's Method

- Randomly choose K data points to be the initial cluster centers



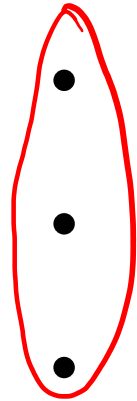
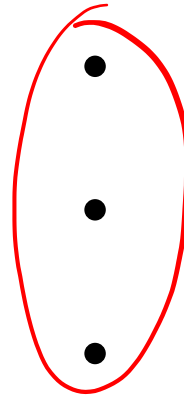
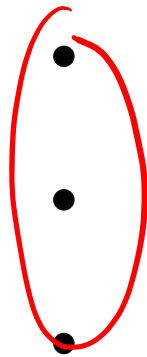
Initializing K -means: Lloyd's Method

- Randomly choose K data points to be the initial cluster centers



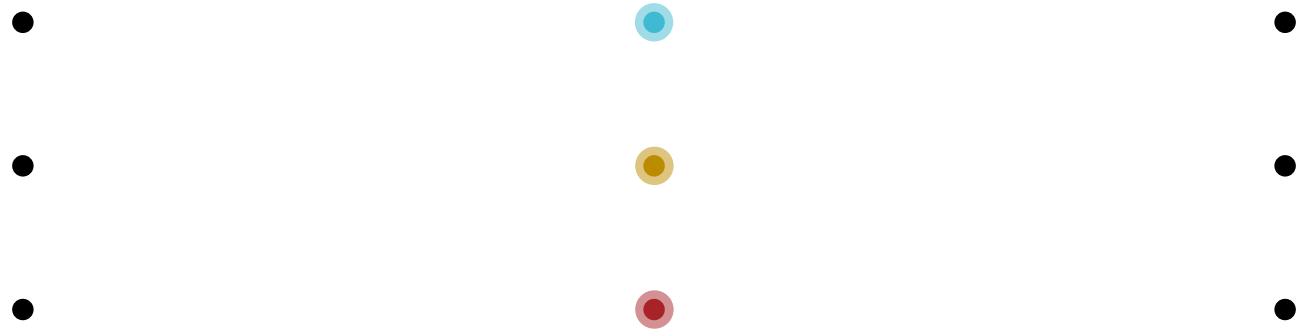
Initializing K -means: Lloyd's Method

- Randomly choose K data points to be the initial cluster centers



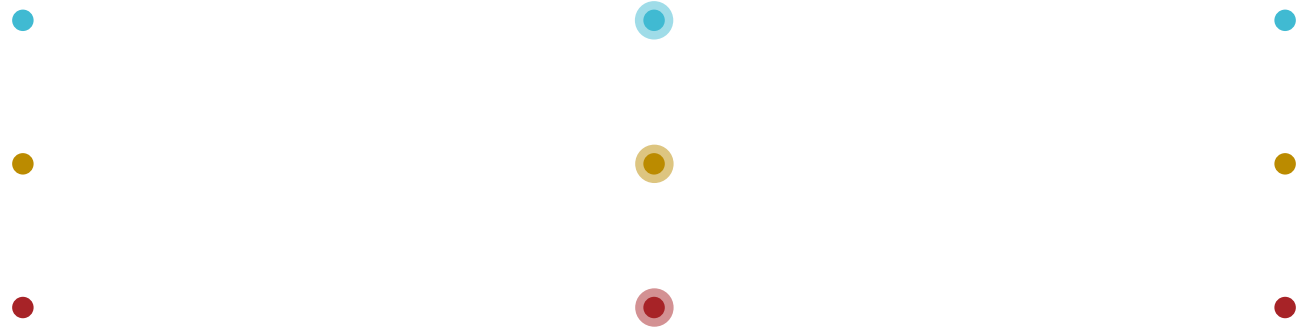
Initializing K -means: Lloyd's Method

- Randomly choose K data points to be the initial cluster centers



Initializing K -means: Lloyd's Method

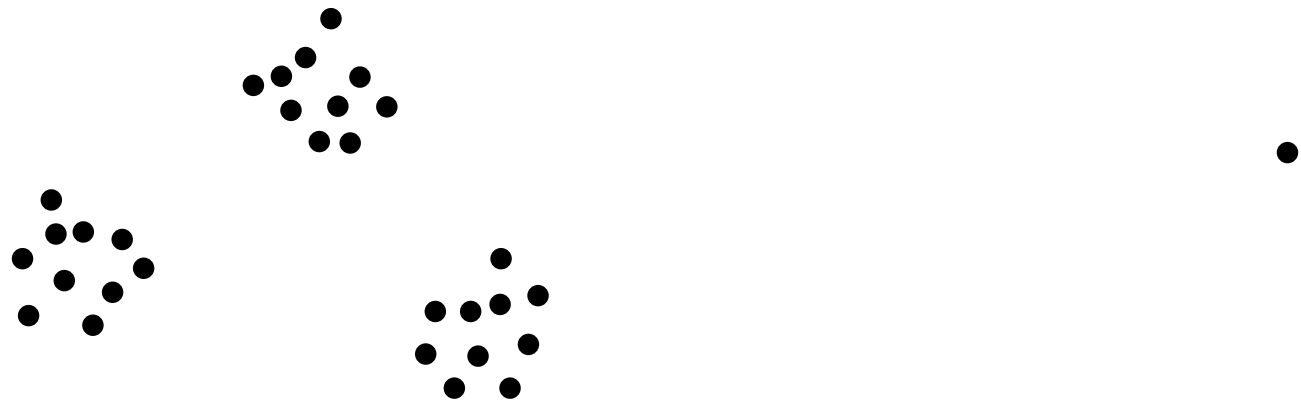
- Randomly choose K data points to be the initial cluster centers



- Lloyd's method converges to a local minimum and that local minimum can be arbitrarily bad (relative to the optimal clusters)
- Intuition: want initial cluster centers to be far apart from one another

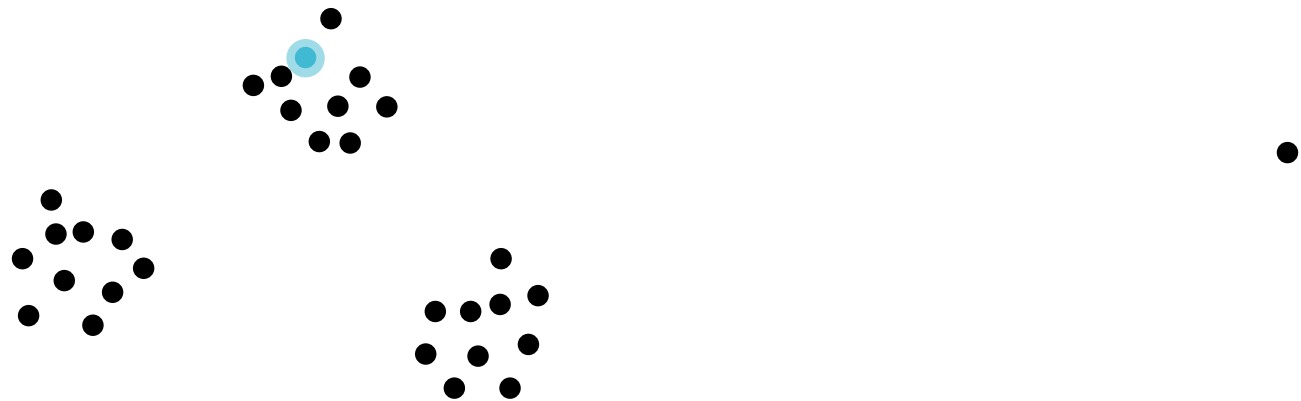
Initializing K -means: Furthest Point

1. Choose the first cluster center randomly from the data points
2. For each other data point \mathbf{x} , compute $D(\mathbf{x})$, the distance between \mathbf{x} and the closest cluster center to \mathbf{x}
3. Select the data point with the largest $D(\mathbf{x})$ as the next cluster center
4. Repeat 2 and 3 $K - 1$ times



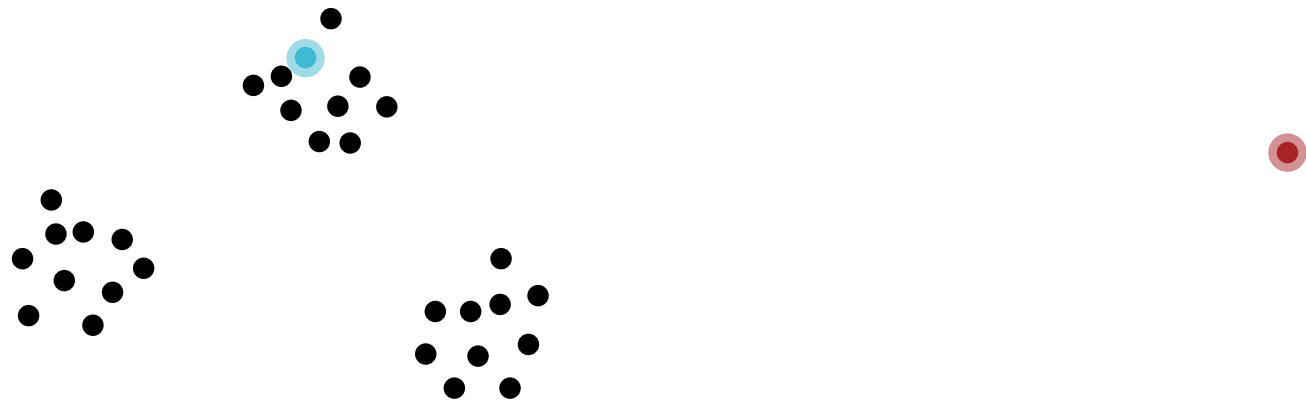
Initializing K -means: Furthest Point

1. Choose the first cluster center randomly from the data points
2. For each other data point x , compute $D(x)$, the distance between x and the closest cluster center to x
3. Select the data point with the largest $D(x)$ as the next cluster center
4. Repeat 2 and 3 $K - 1$ times



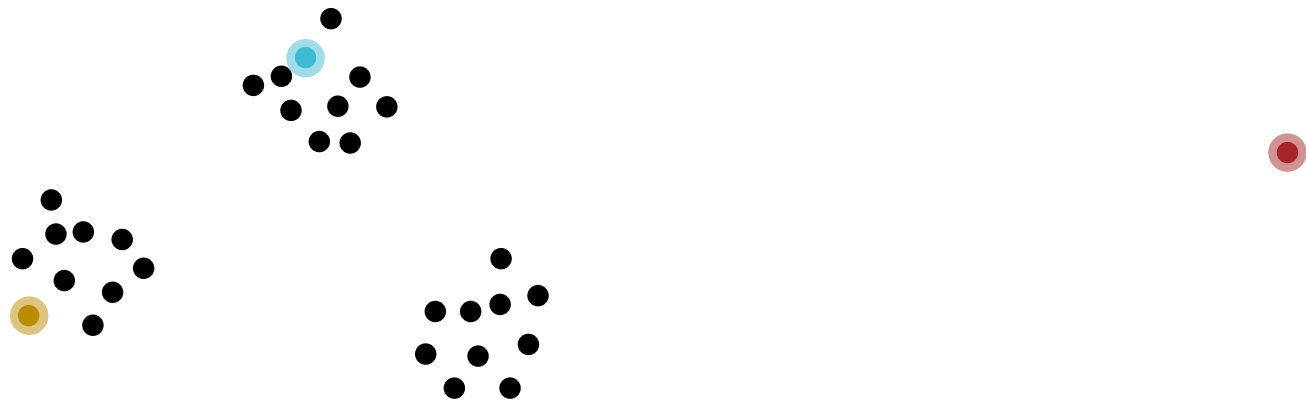
Initializing K -means: Furthest Point

1. Choose the first cluster center randomly from the data points
2. For each other data point x , compute $D(x)$, the distance between x and the closest cluster center to x
3. Select the data point with the largest $D(x)$ as the next cluster center
4. Repeat 2 and 3 $K - 1$ times



Initializing K -means: Furthest Point

1. Choose the first cluster center randomly from the data points
2. For each other data point x , compute $D(x)$, the distance between x and the closest cluster center to x
3. Select the data point with the largest $D(x)$ as the next cluster center
4. Repeat 2 and 3 $K - 1$ times



Initializing K -means: Furthest Point

1. Choose the first cluster center randomly from the data points
2. For each other data point \mathbf{x} , compute $D(\mathbf{x})$, the distance between \mathbf{x} and the closest cluster center to \mathbf{x}
3. Select the data point with the largest $D(\mathbf{x})$ as the next cluster center
4. Repeat 2 and 3 $K - 1$ times



Initializing K -means: Furthest Point

1. Choose the first cluster center randomly from the data points
2. For each other data point \mathbf{x} , compute $D(\mathbf{x})$, the distance between \mathbf{x} and the closest cluster center to \mathbf{x}
3. Select the data point with the largest $D(\mathbf{x})$ as the next cluster center
4. Repeat 2 and 3 $K - 1$ times
 - Works great in the case of well-clustered data!
 - Can struggle with outliers...

Initializing K -means: K -means++

1. Choose the first cluster center randomly from the data points
2. For each other data point \mathbf{x} , compute $D(\mathbf{x})$, the distance between \mathbf{x} and the closest cluster center to \mathbf{x}
3. *Sample* the next cluster center with probability proportional to $D(\mathbf{x})^2$
4. Repeat 2 and 3 $K - 1$ times

i	$D(\mathbf{x}^{(i)})$	$D(\mathbf{x}^{(i)})^2$	Probability of Being Selected
1	4	16	16/123
2	7	49	49/123
\vdots	\vdots	\vdots	\vdots
N	1	1	1/123
Total		123	123/123 = 1

Initializing K -means: K -means++

1. Choose the first cluster center randomly from the data points
 2. For each other data point \mathbf{x} , compute $D(\mathbf{x})$, the distance between \mathbf{x} and the closest cluster center to \mathbf{x}
 3. *Sample* the next cluster center with probability proportional to $D(\mathbf{x})^2$
 4. Repeat 2 and 3 $K - 1$ times
- K -means++ achieves a $O(\log K)$ approximation to the optimal clustering in expectation!
 - All initialization methods can benefit from multiple random restarts

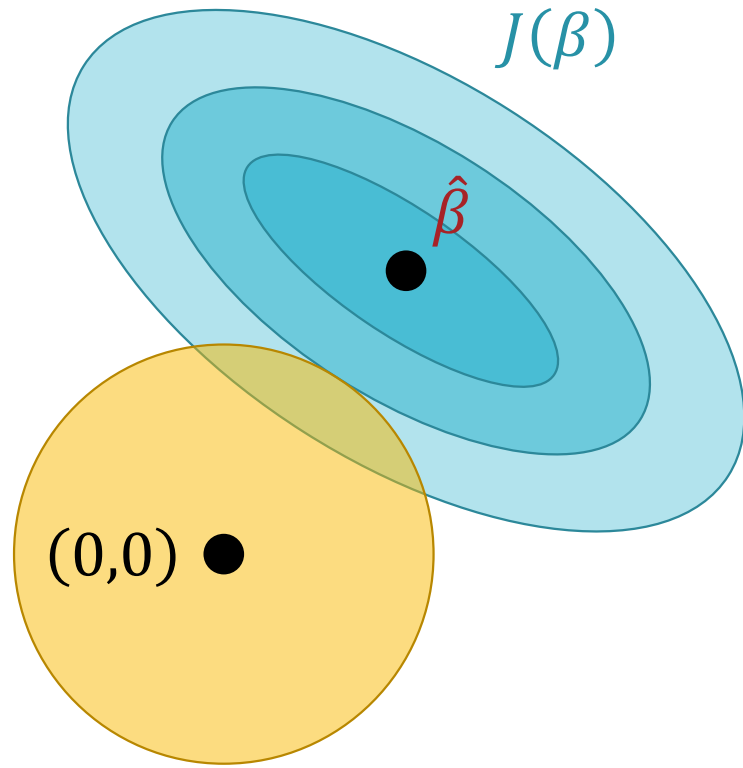
K-means Learning Objectives

You should be able to...

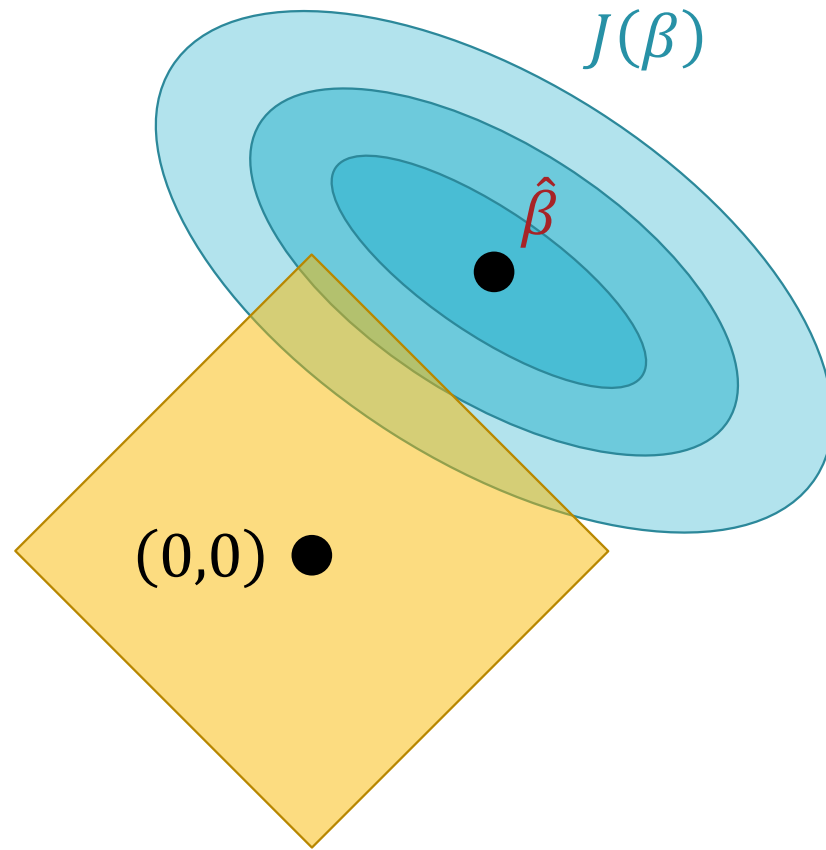
- Distinguish between coordinate descent and block coordinate descent
- Define an objective function that gives rise to a "good" clustering
- Apply block coordinate descent to an objective function preferring each point to be close to its nearest objective function to obtain the K-Means algorithm
- Implement the K-Means algorithm
- Connect the nonconvexity of the K-Means objective function with the (possibly) poor performance of random initialization

Dimensionality Reduction

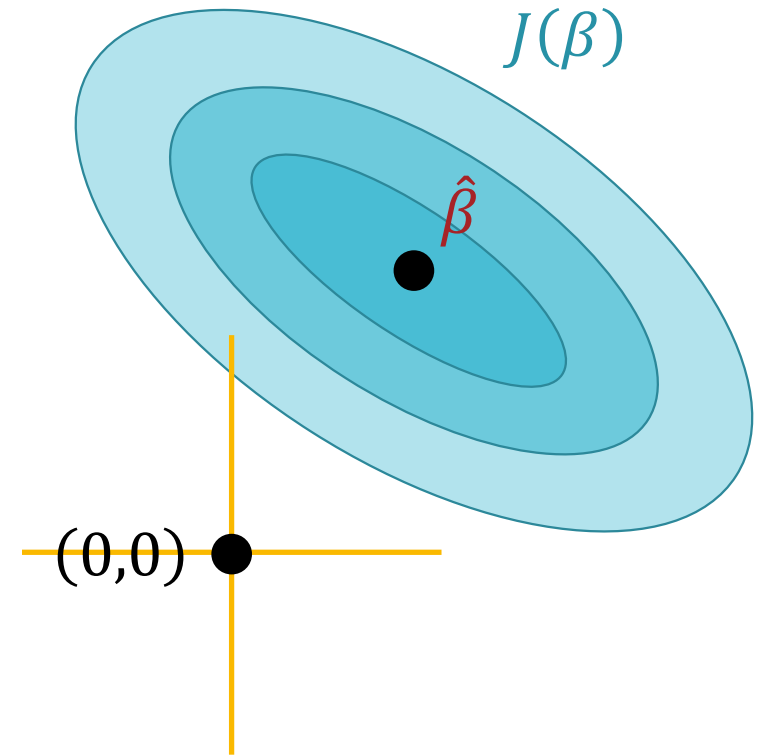
- Goal: given some unlabeled data set, learn a latent (typically lower-dimensional) representation
- Use cases:
 - Reducing computational cost (runtime, storage, etc...)
 - Improving generalization
 - Visualizing data
- Applications:
 - word embeddings
 - image/video/audio
 - ratings on social media platforms



Ridge or L_2

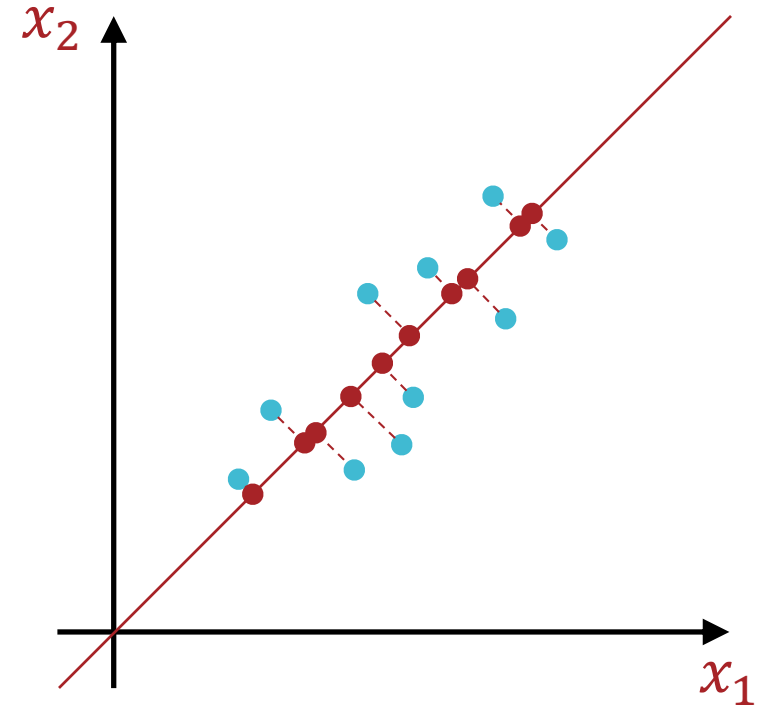
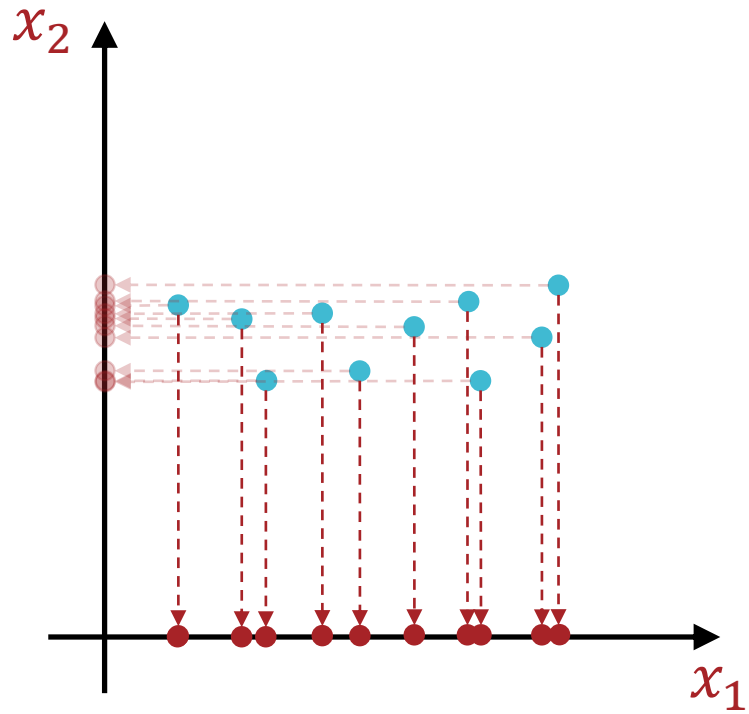


Lasso or L_1

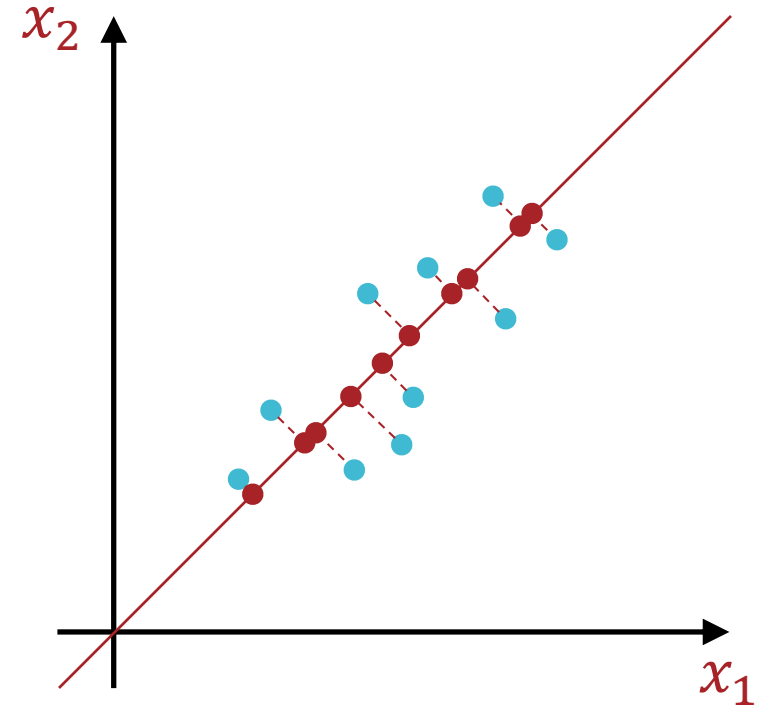
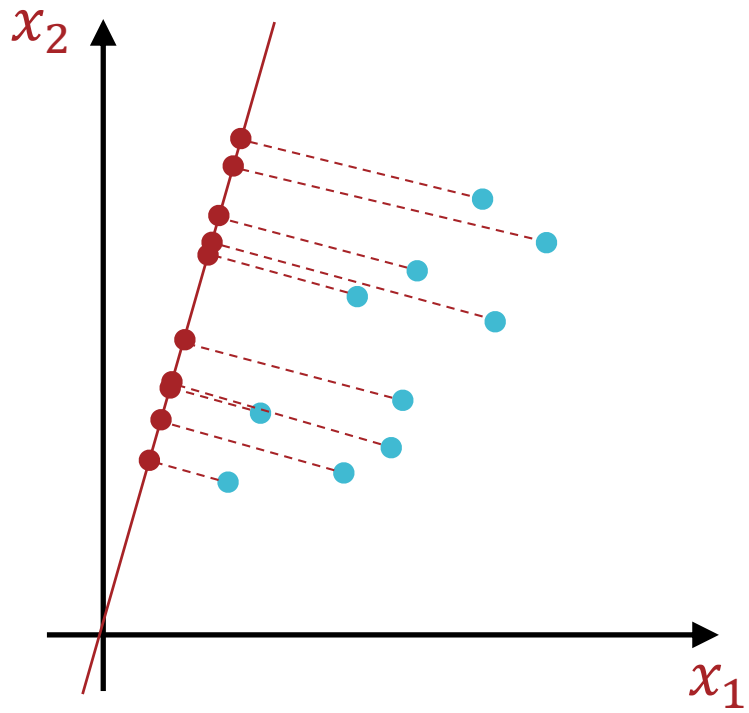


L_0

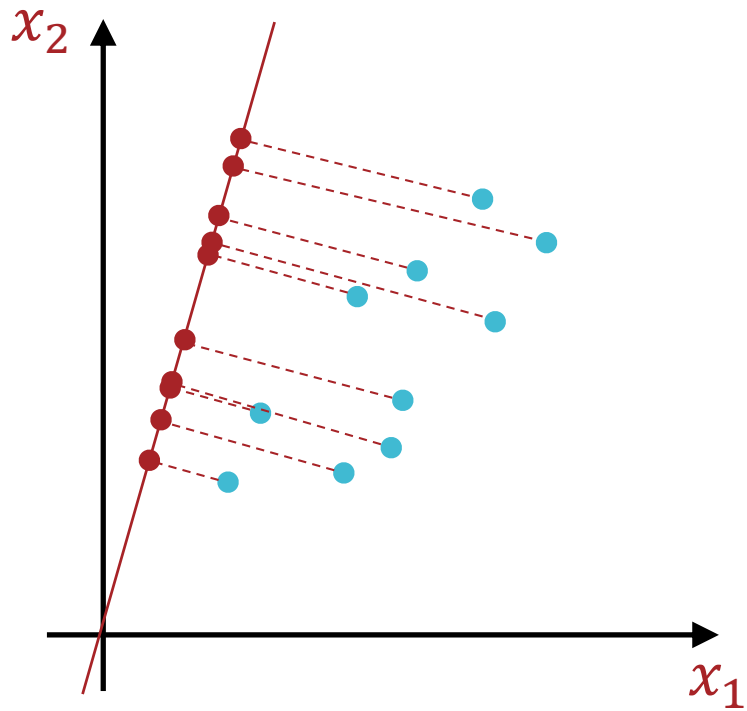
Recall: L_1 (or L_0) Regularization



Feature Elimination

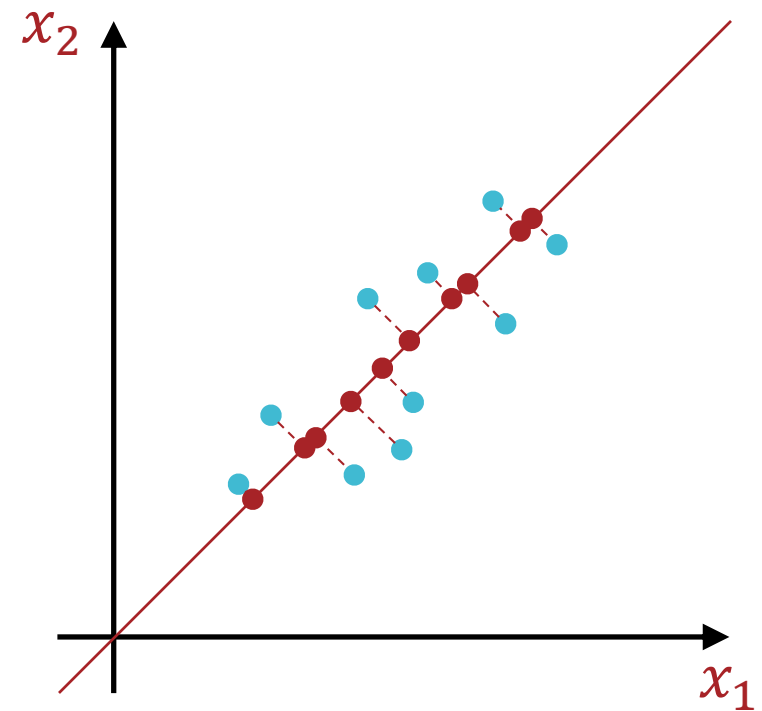


Feature Reduction



Option A

Option B (TOXIC)



Option C

Which projection do you prefer (Q1) and why (Q2)?

Background: Sample Variance and Covariance

- Given a collection of N 1-dimensional samples $[x^{(1)}, x^{(2)}, \dots, x^{(N)}]$ from some random variable, the **sample variance** is

$$\rightarrow \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N \underbrace{(x^{(i)} - \hat{\mu})^2}_{\text{red bracket}} = \frac{1}{N} \sum_{i=1}^N \left(x^{(i)} - \frac{1}{N} \sum_{n=1}^N x^{(n)} \right)^2$$

- Given a collection of N D -dimensional samples $[\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}]$ from some random variable, the **sample covariance** between dimension j and k is

$$\hat{\Sigma}_{jk} = \frac{1}{N} \sum_{i=1}^N \underbrace{(x_j^{(i)} - \hat{\mu}_j)}_{\text{red bracket}} (x_k^{(i)} - \hat{\mu}_k) \text{ where } \hat{\mu}_d = \frac{1}{N} \sum_{n=1}^N x_d^{(n)}$$

Background: Sample Variance and Covariance

- Given a collection of N 1-dimensional samples $[x^{(1)}, x^{(2)}, \dots, x^{(N)}]$ from some random variable, the **sample variance** is

$$\rightarrow \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \hat{\mu})^2 = \frac{1}{N} \sum_{i=1}^N \left(x^{(i)} - \frac{1}{N} \sum_{n=1}^N x^{(n)} \right)^2$$

- Given a collection of N D -dimensional samples $[\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}]$ from some random variable, the **sample covariance matrix** is

$$\rightarrow \Sigma = \frac{1}{N} X^T X \quad \text{where } X = \begin{matrix} \rightarrow \\ \left[\begin{array}{c} (\mathbf{x}^{(1)} - \boldsymbol{\mu})^T \\ (\mathbf{x}^{(2)} - \boldsymbol{\mu})^T \\ \vdots \\ (\mathbf{x}^{(N)} - \boldsymbol{\mu})^T \end{array} \right] \end{matrix}$$

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

$$(a-b)^T = a^T - b^T$$

Reconstruction Error

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

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

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

$$\begin{aligned} & (\tilde{\mathbf{x}}^{(n)} - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v})^T (\tilde{\mathbf{x}}^{(n)} - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}) \\ \Rightarrow & \tilde{\mathbf{x}}^{(n)T} \tilde{\mathbf{x}}^{(n)} - 2(\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}^T \tilde{\mathbf{x}}^{(n)} + (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)})^2 \mathbf{v}^T \mathbf{v} \\ = & \tilde{\mathbf{x}}^{(n)T} \tilde{\mathbf{x}}^{(n)} - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)})^2 \end{aligned}$$

Minimizing the Reconstruction Error

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

$$\Rightarrow \operatorname{argmin}_{v: \|v\|_2=1} \sum_{n=1}^N \underbrace{\tilde{x}^{(n)T} \tilde{x}^{(n)}}_{=0} - (v^T \tilde{x}^{(n)})^2$$

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

is equivalent to maximizing the variance of the projections

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

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

$X^T X$

Maximizing the Variance

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

$$\star L(v) = v^T X^T X v - \lambda (v^T v - 1)$$

$$\frac{\partial L}{\partial v} = 2X^T X v - 2\lambda v$$

$$\Rightarrow 2X^T X \hat{v} - 2\lambda \hat{v} = 0$$

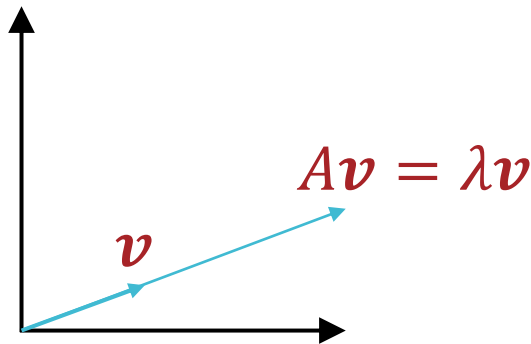
$$\Rightarrow X^T X \hat{v} = \lambda \hat{v}$$

\hat{v} is an eigenvector for the matrix $X^T X$ w/ eigenvalue λ

Background: Eigenvectors & Eigenvalues

- Given a square matrix $A \in \mathbb{R}^{N \times N}$, a vector $\mathbf{v} \in \mathbb{R}^{N \times 1}$ is an **eigenvector** of A iff there exists some scalar λ such that

$$A\mathbf{v} = \lambda\mathbf{v}$$



Intuition: A scales or stretches \mathbf{v} but does not rotate it

- Key property: the eigenvectors of symmetric matrices (e.g., the covariance matrix of a data set) are orthogonal!

Maximizing the Variance

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

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

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}} = 2(X^T X) \mathbf{v} - 2\lambda \mathbf{v}$$

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

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

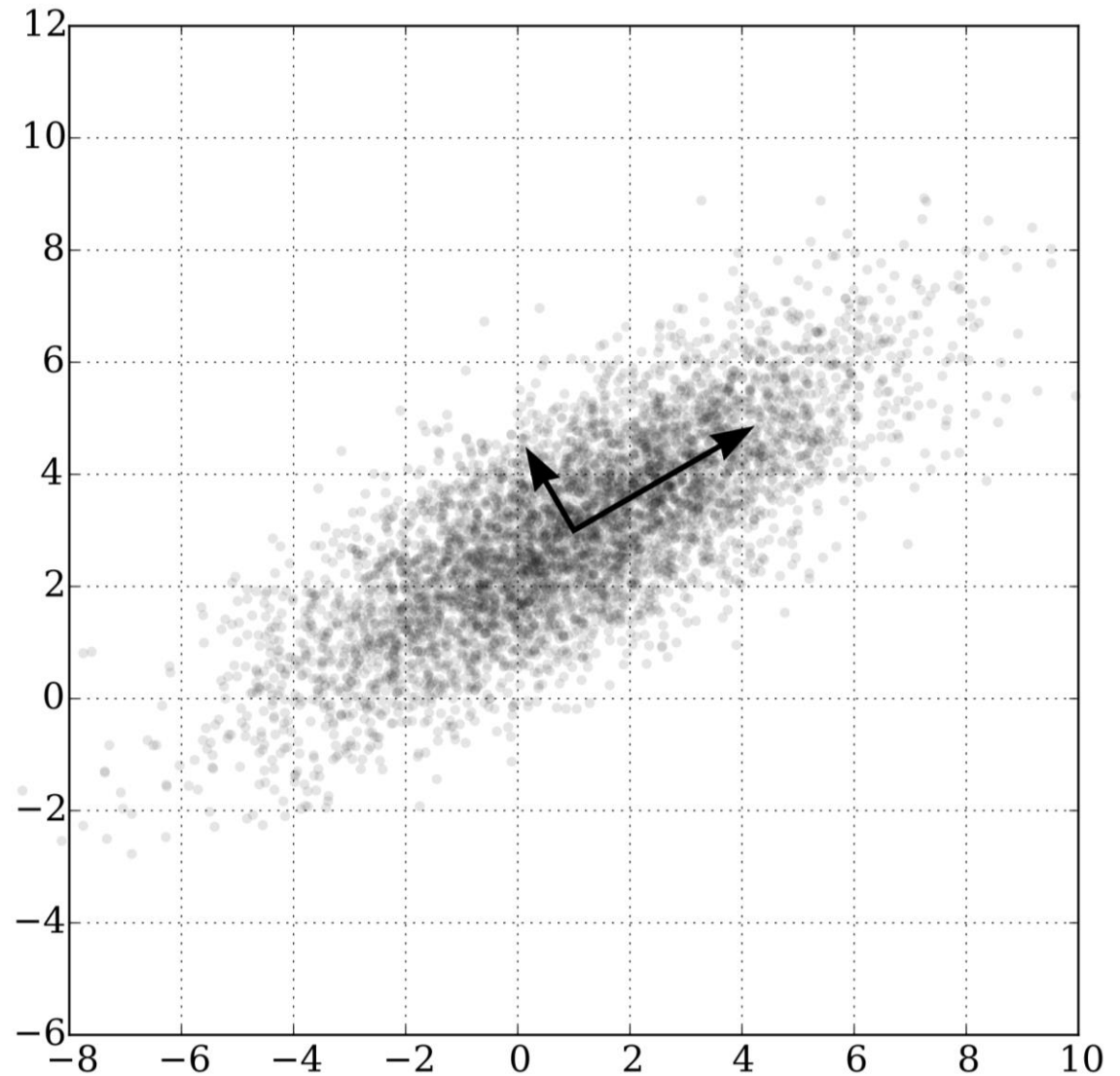
Maximizing the Variance

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

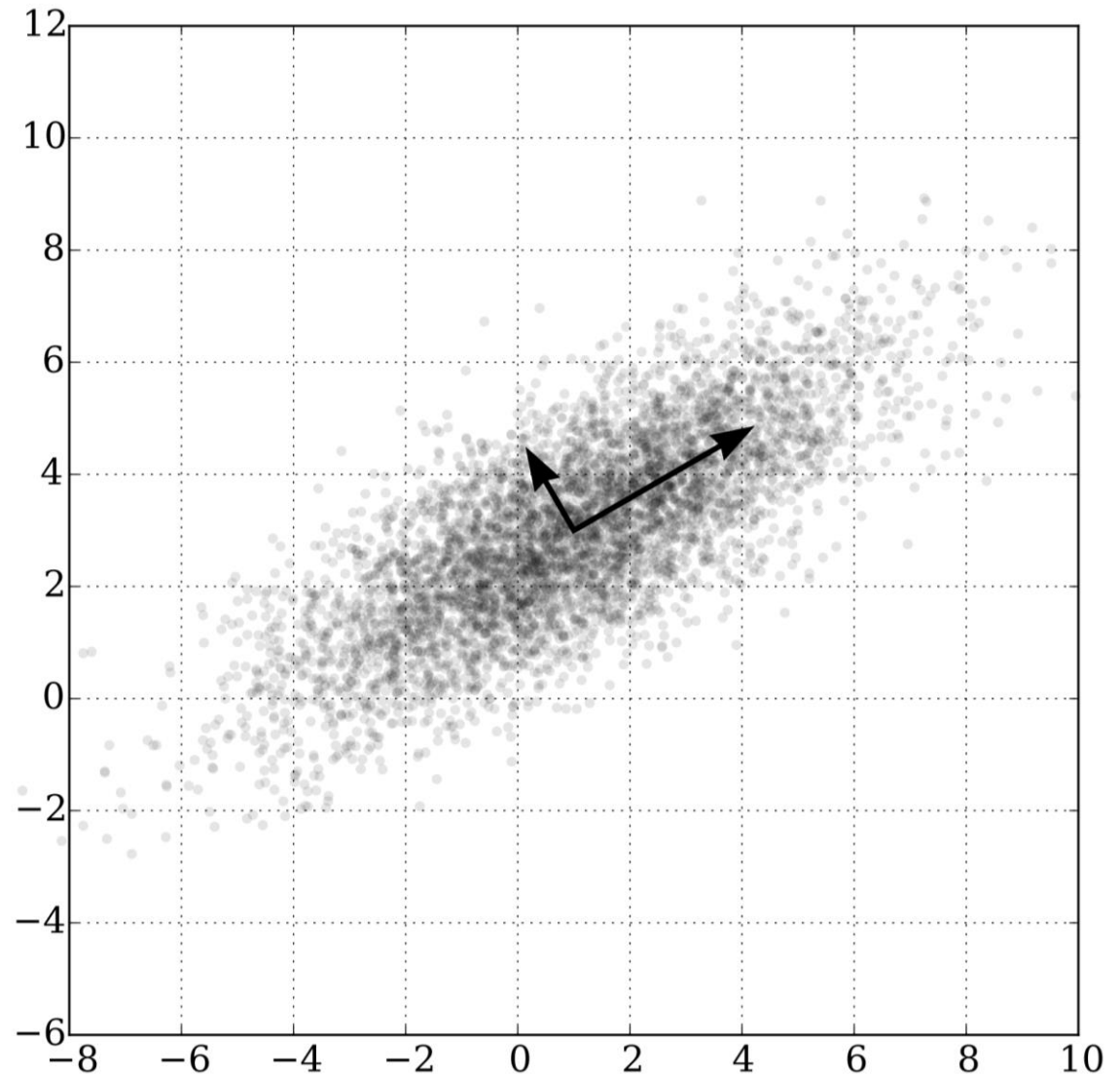
$$(X^T X) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}} \rightarrow \hat{\mathbf{v}}^T (X^T X) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}}^T \hat{\mathbf{v}} = \lambda$$

- The first principal component is the eigenvector $\hat{\mathbf{v}}_1$ that corresponds to the largest eigenvalue 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 , $Z = XV_\rho$
- Output: Z , the transformed (potentially lower-dimensional) data

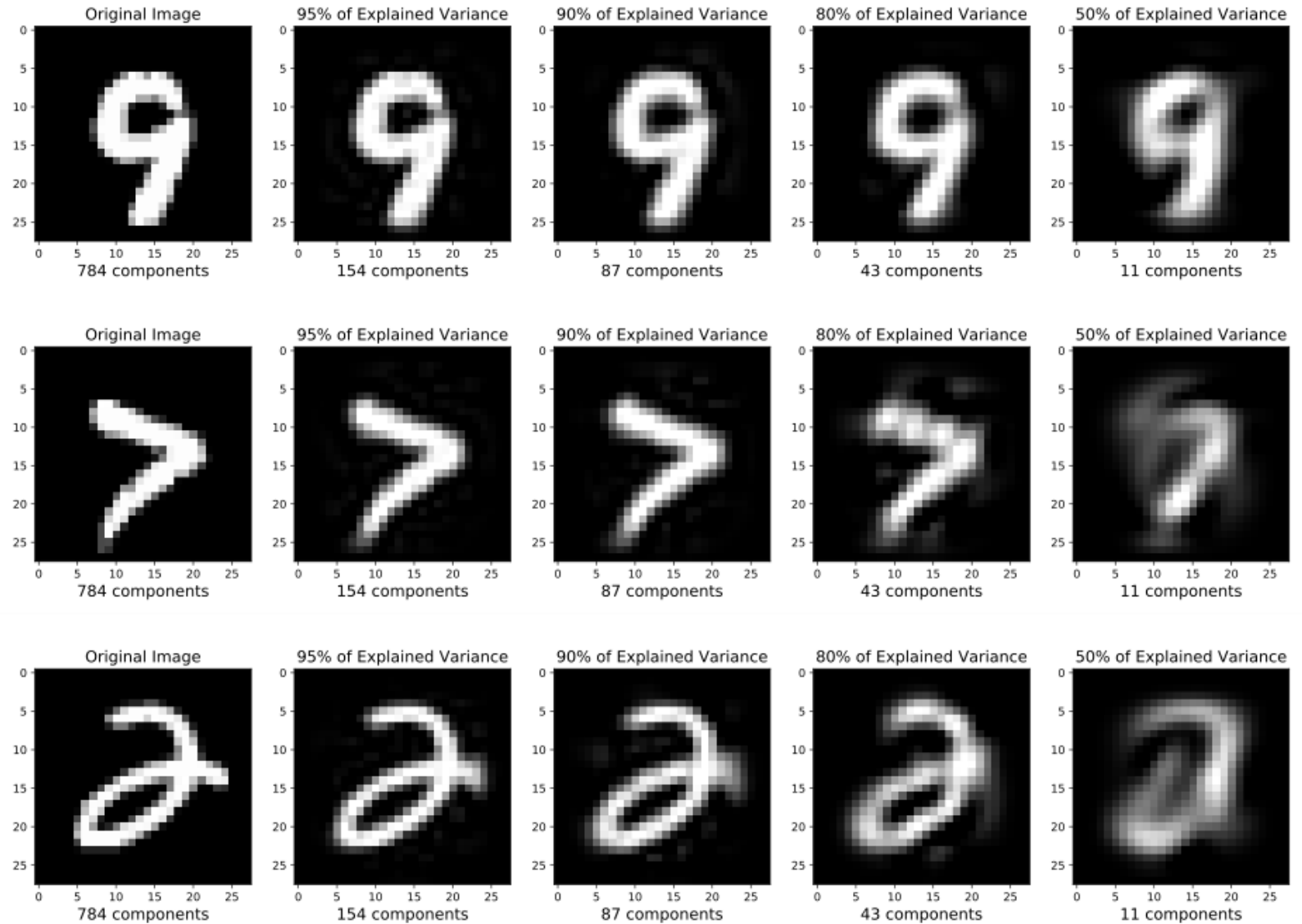
Choosing the number of PCs

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

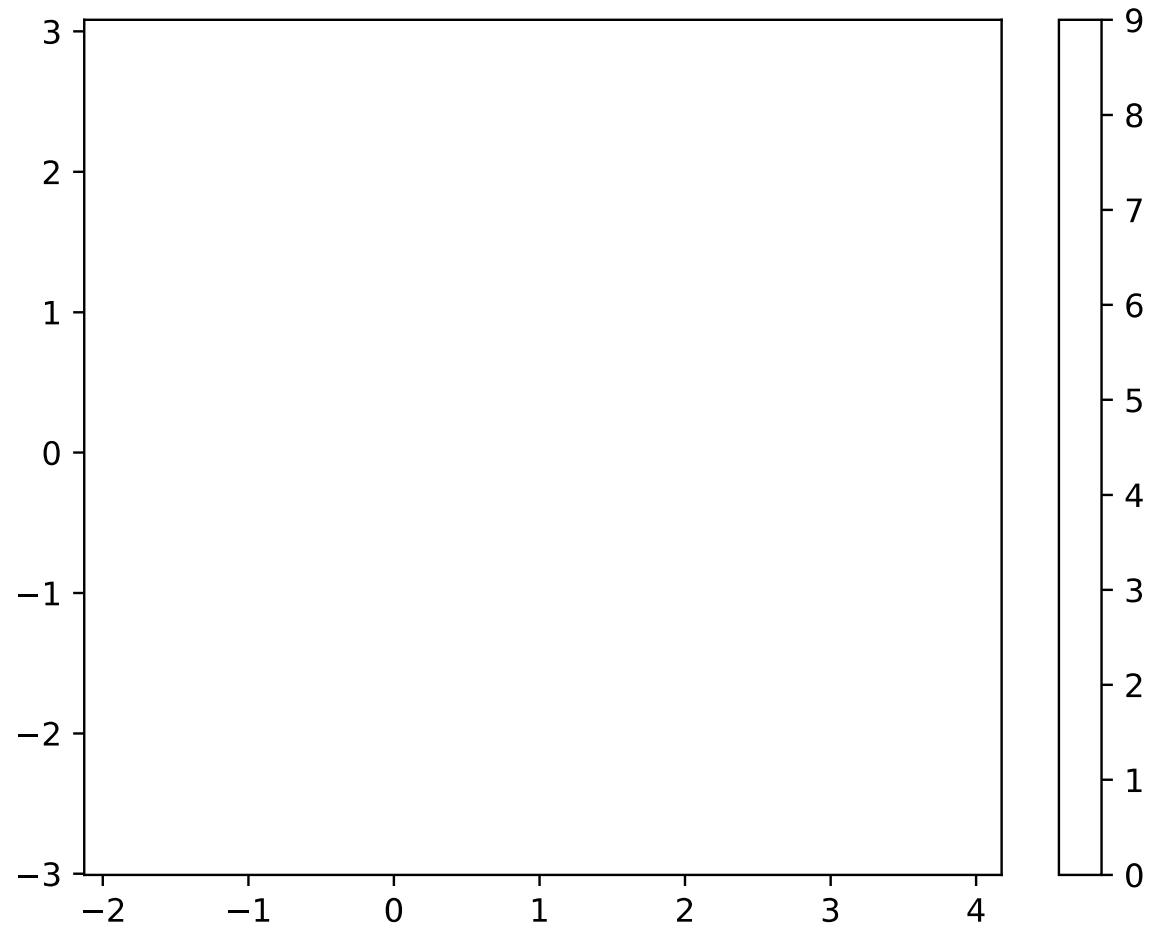
$$\lambda_i / \sum \lambda_j$$

- 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

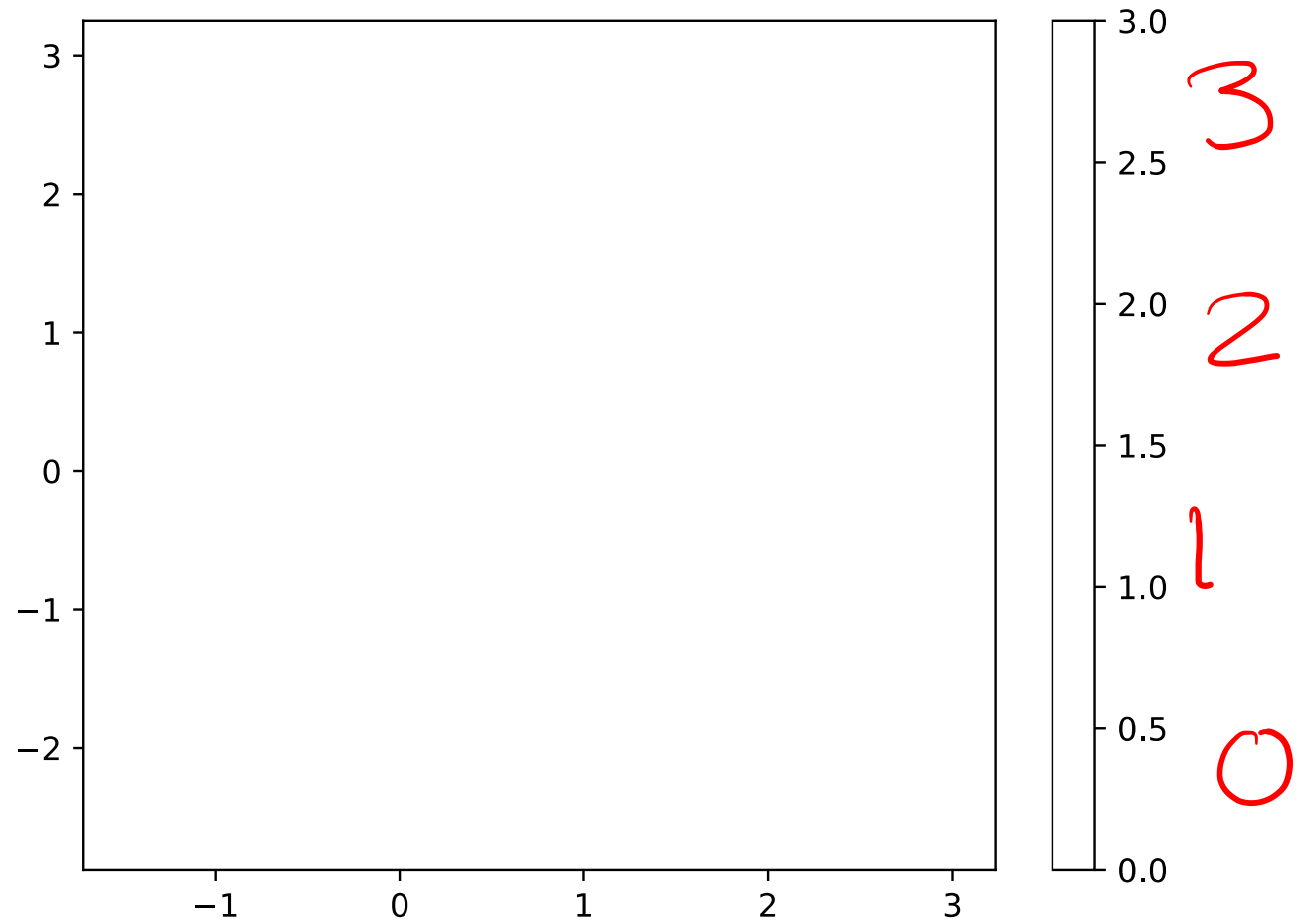
PCA Example: MNIST Digits



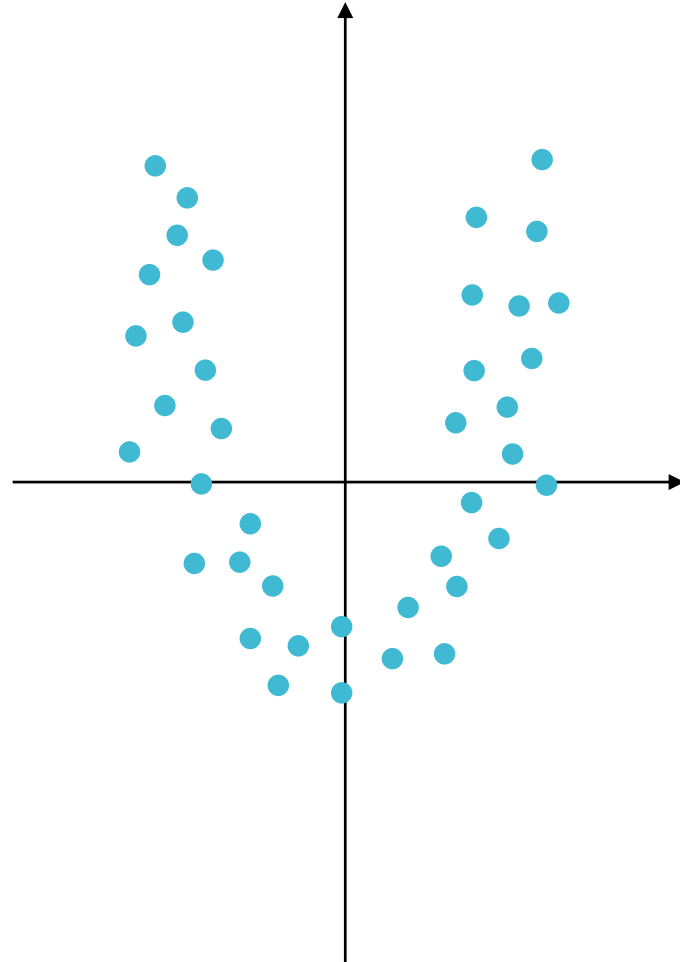
PCA Example: MNIST Digits



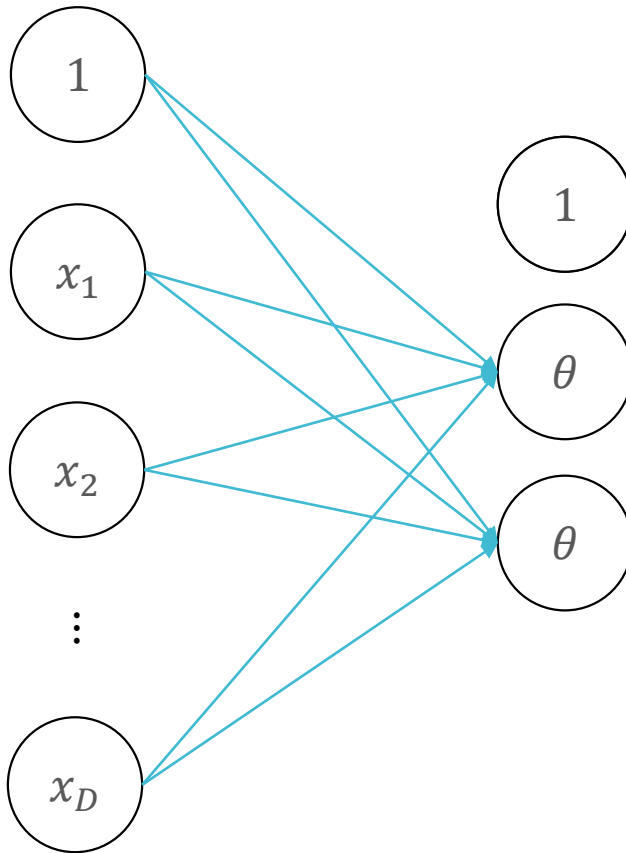
PCA Example: MNIST Digits



Shortcomings of PCA

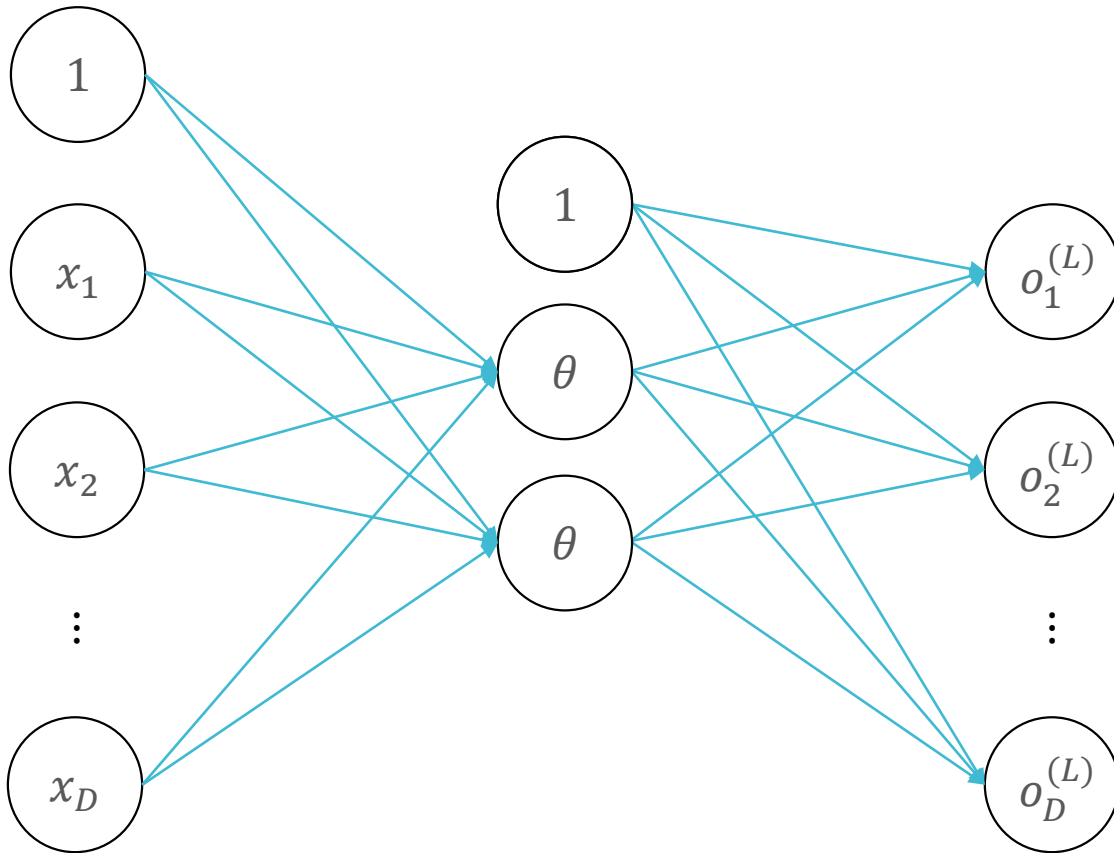


Autoencoders



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

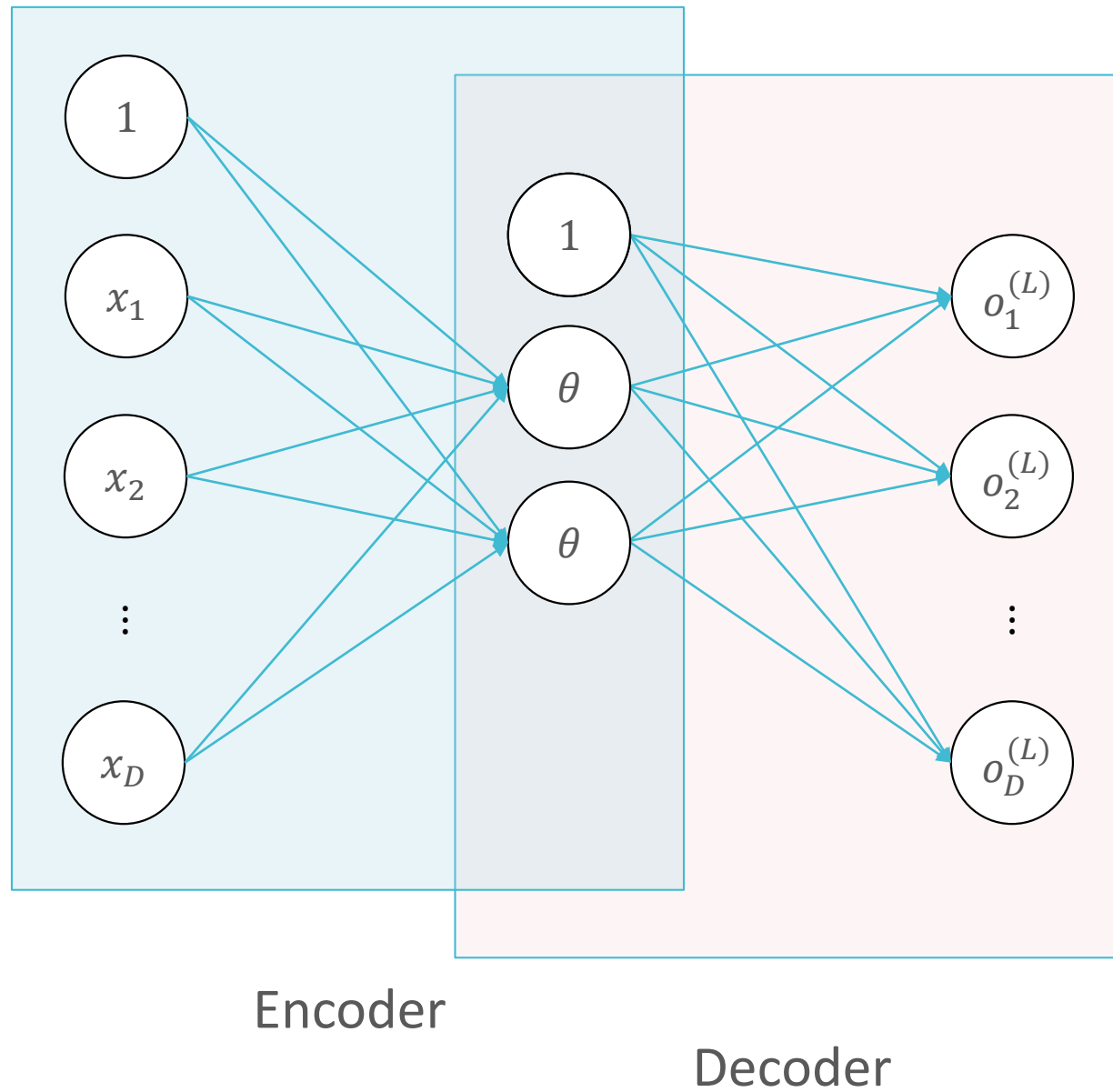
Autoencoders



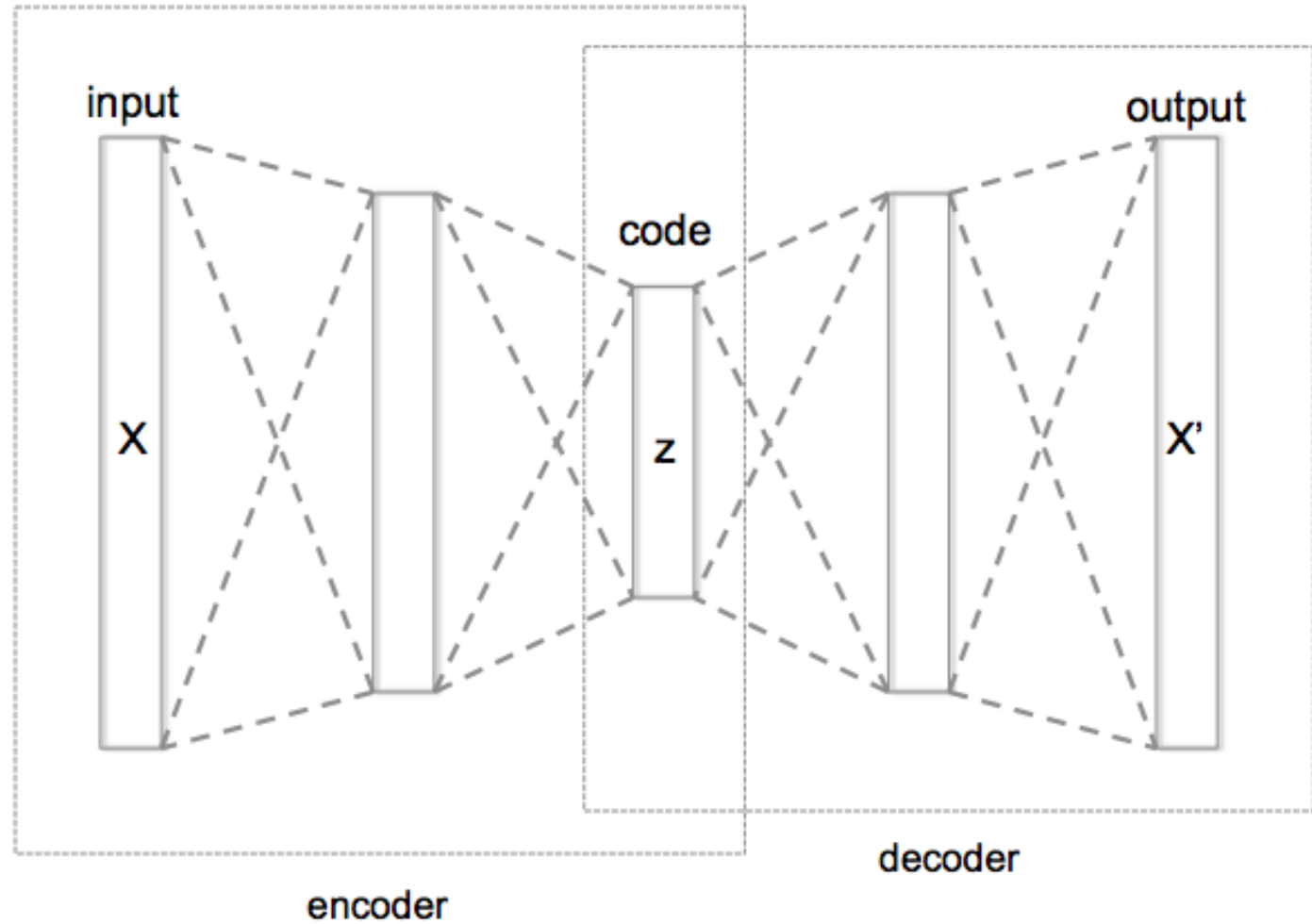
- Learn the weights by minimizing the reconstruction loss:

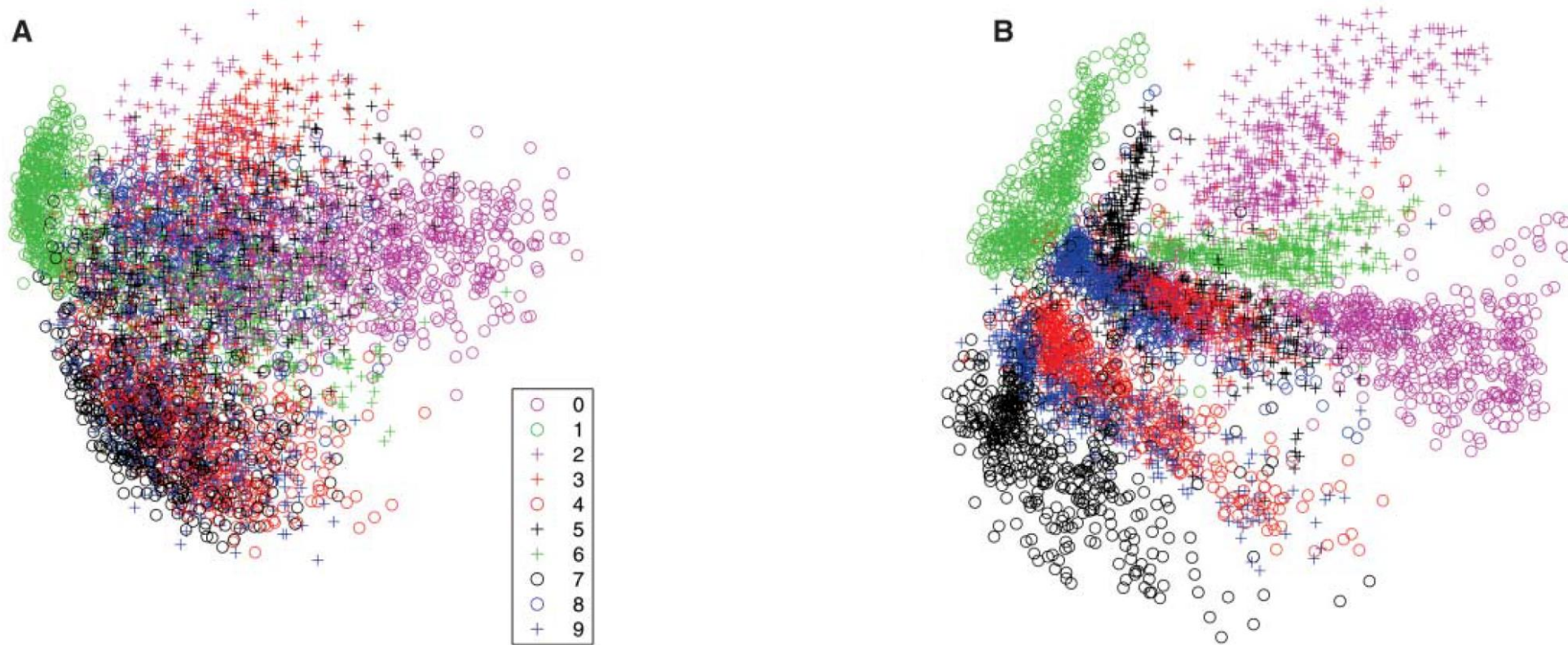
$$e(\mathbf{x}) = \|\mathbf{x} - \mathbf{o}^{(L)}\|_2^2$$

Autoencoders



Deep Autoencoders





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

PCA Learning Objectives

You should be able to...

- Define the sample mean, sample variance, and sample covariance of a vector-valued dataset
- Identify examples of high dimensional data and common use cases for dimensionality reduction
- Draw the principal components of a given toy dataset
- Establish the equivalence of minimization of reconstruction error with maximization of variance
- Given a set of principal components, project from high to low dimensional space and do the reverse to produce a reconstruction
- Explain the connection between PCA, eigenvectors, eigenvalues, and covariance matrix
- Use common methods in linear algebra to obtain the principal components