

Dimensionality Reduction

Lecture 12

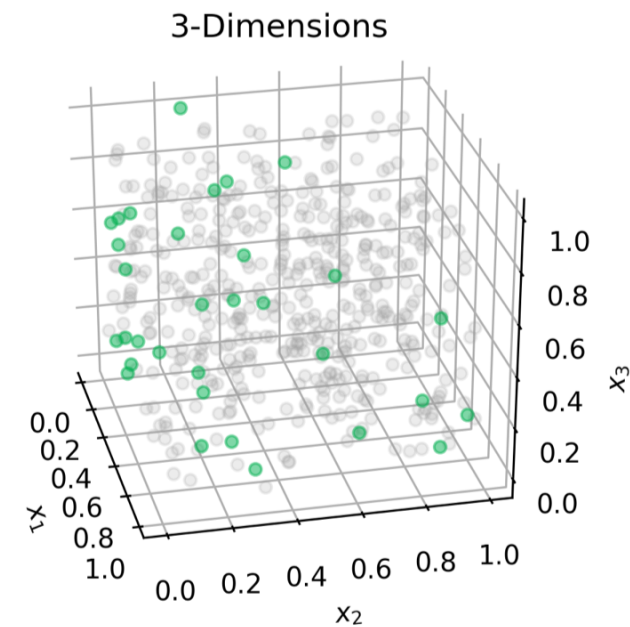
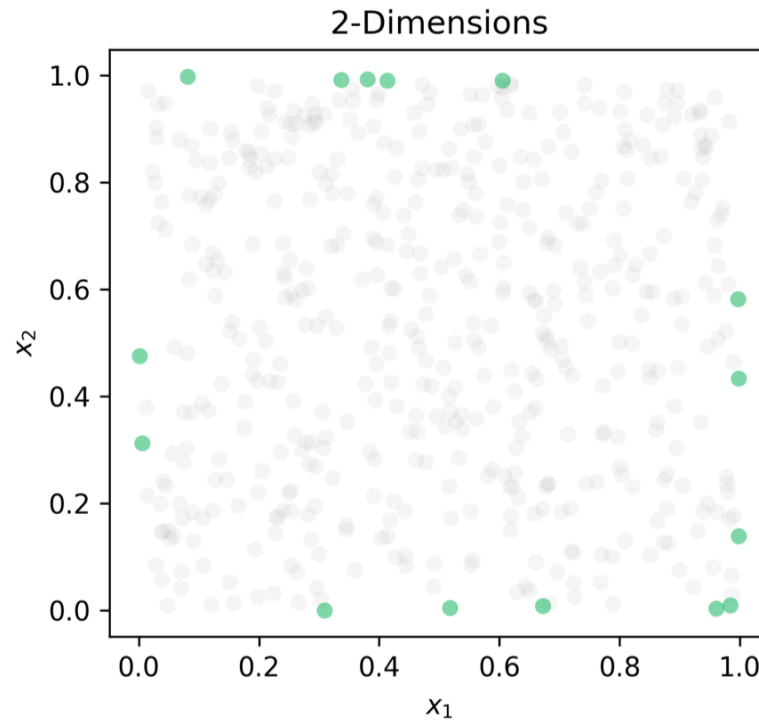
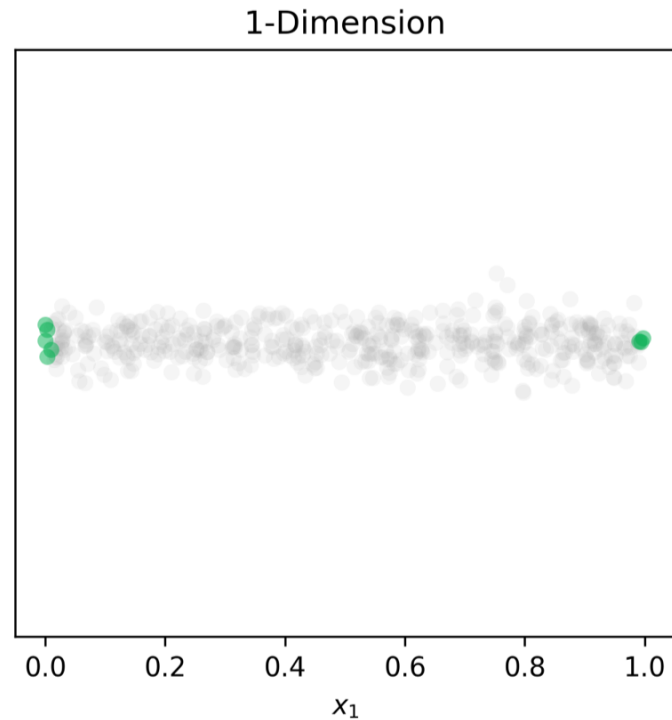
The **Curse** of Dimensionality

Challenge 1

In high dimensions, data become sparse
(increasing the risk of overfitting)

Random data points in a unit hypercube...

- Data point is a distance < 0.01 units from the edge of a unit hypercube
- All other data



Fraction
of edge
data

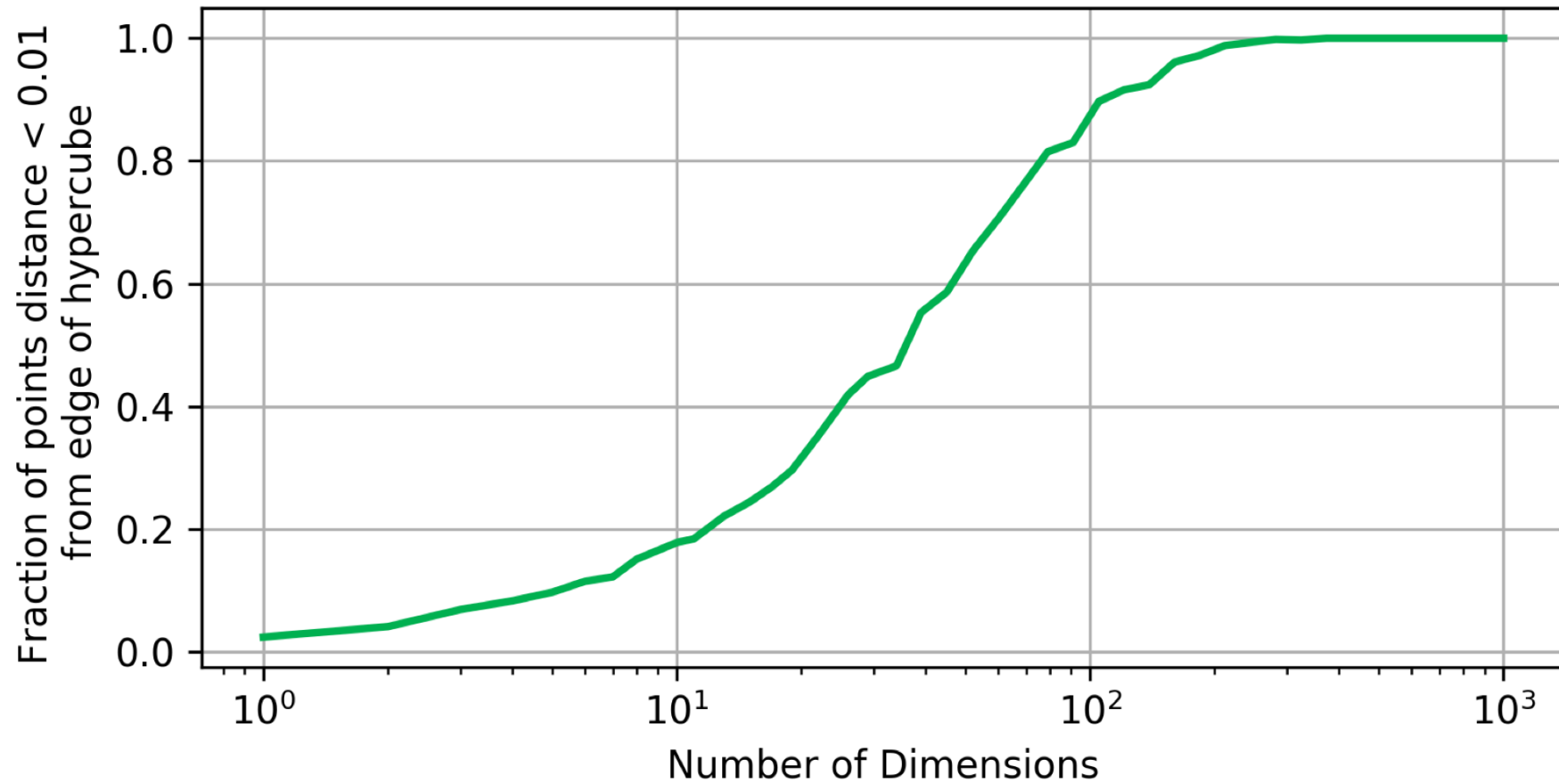
● + ● =

0.016

0.030

0.064

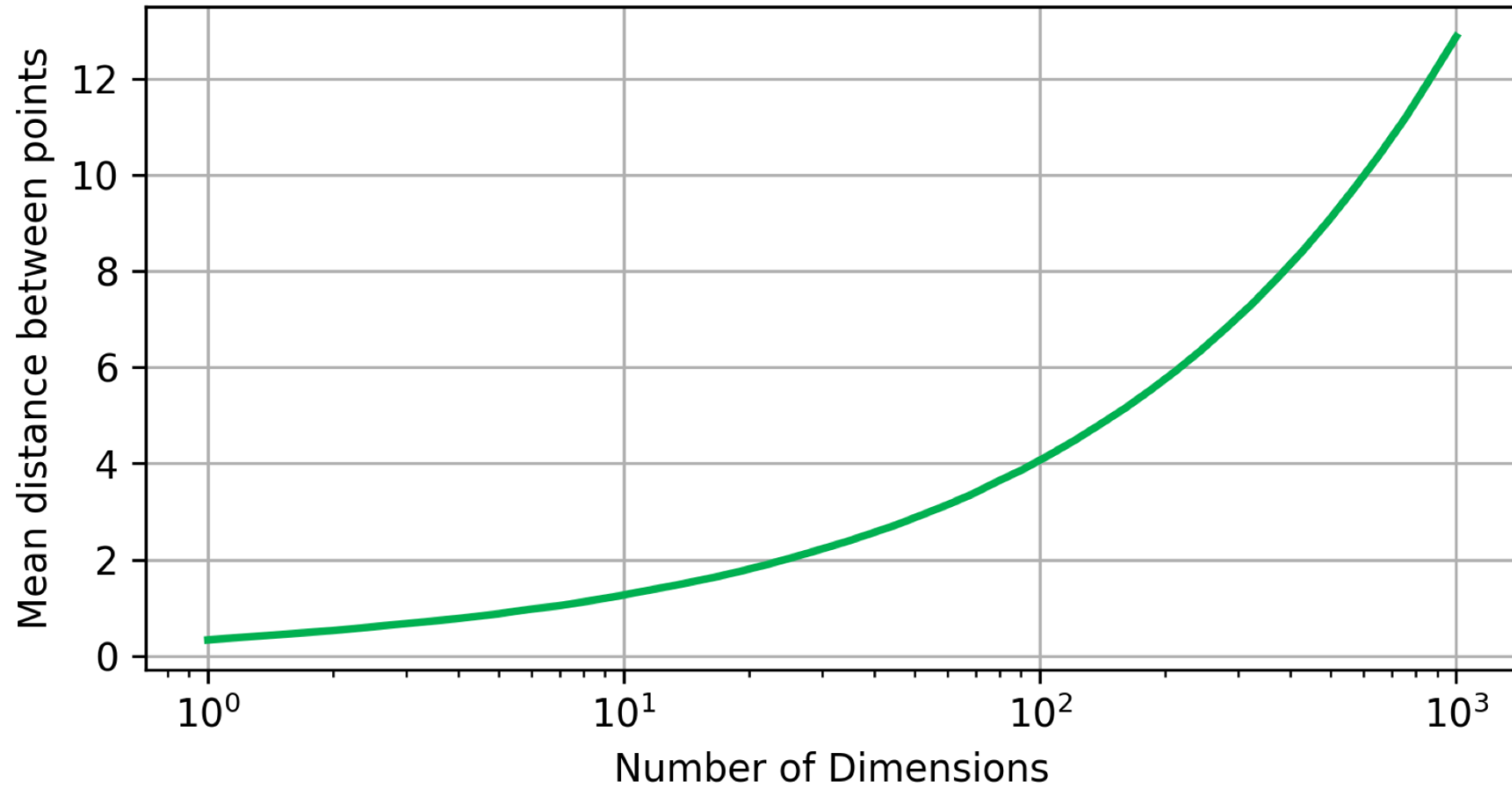
In high dimensions...



...nearly all of the high dimensional space is **far away from the center**

Note: figures constructed using 1,000 random points

In high dimensions...



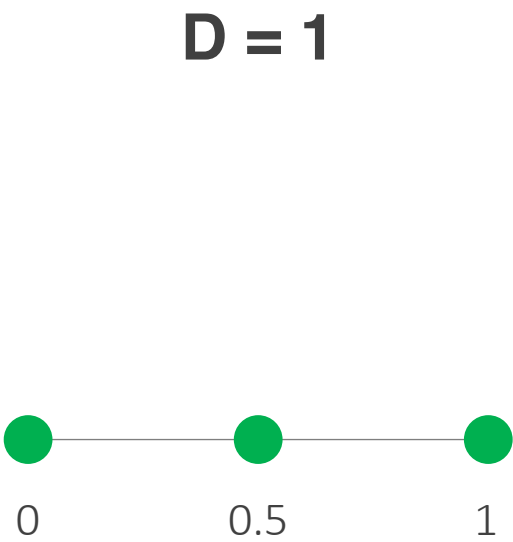
...data become sparse

Note: figures constructed using 1,000 random points

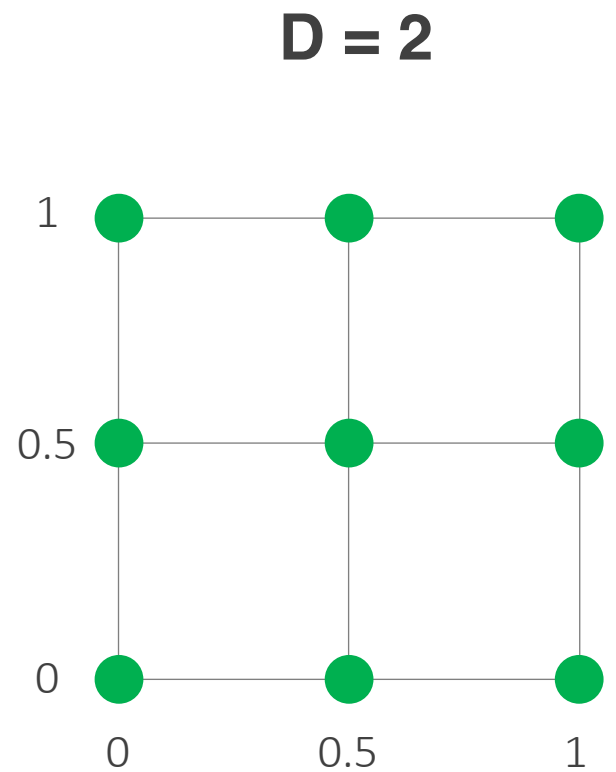
Challenge 2

Much more data are needed for sampling higher dimensional spaces

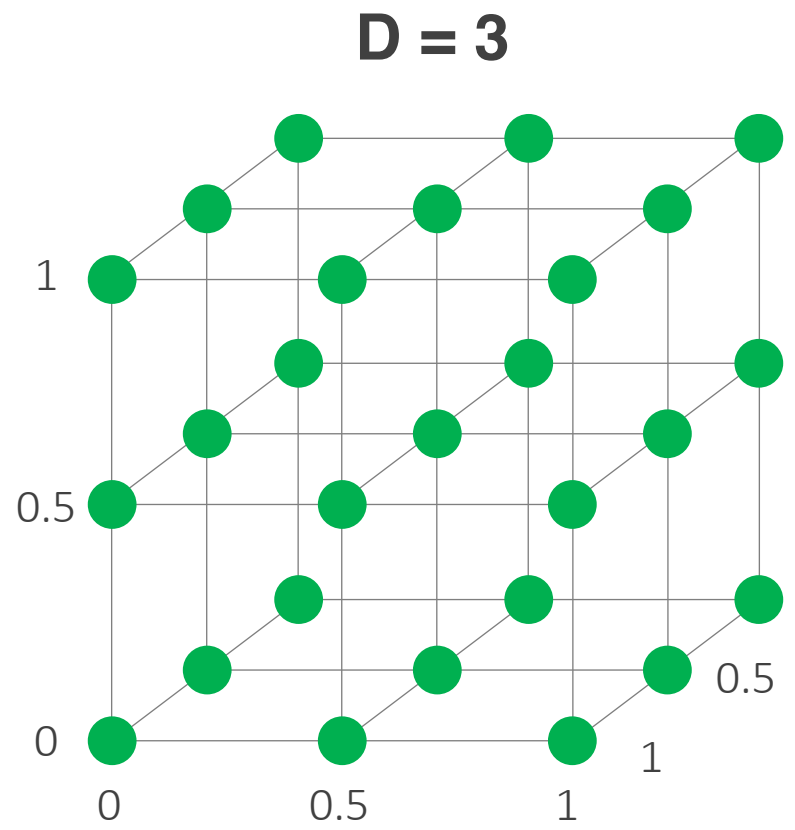
Sample a unit hypercube on a grid spaced at intervals of 0.5



N = 3



N = 9



N = 27

Dim (D)	Samples (N)
4	81
5	243
10	59,049
100	5.2×10^{47}
300	1.4×10^{143}

1 million Googles,
each with 20 exabytes
of data would only be
 10^{25} bytes

...it takes more data to learn in high dimensional spaces

Dimensionality Reduction

Benefits:

- Simplified data processing
- Reduced redundancy of features
- Improved numerical stability due to removed correlations

Approaches:

- Feature selection (subset selection)
- Feature extraction (including some regularization techniques)

Popular approach:

- Principal Components Analysis (PCA)

PCA

Before you begin: Normalize the data!

For each feature, subtract the mean and divide by the standard deviation

Normalized feature

Raw, unscaled feature

$$\mathcal{X} = \frac{x' - \bar{x'}}{\sigma_{x'}}$$

columns = features

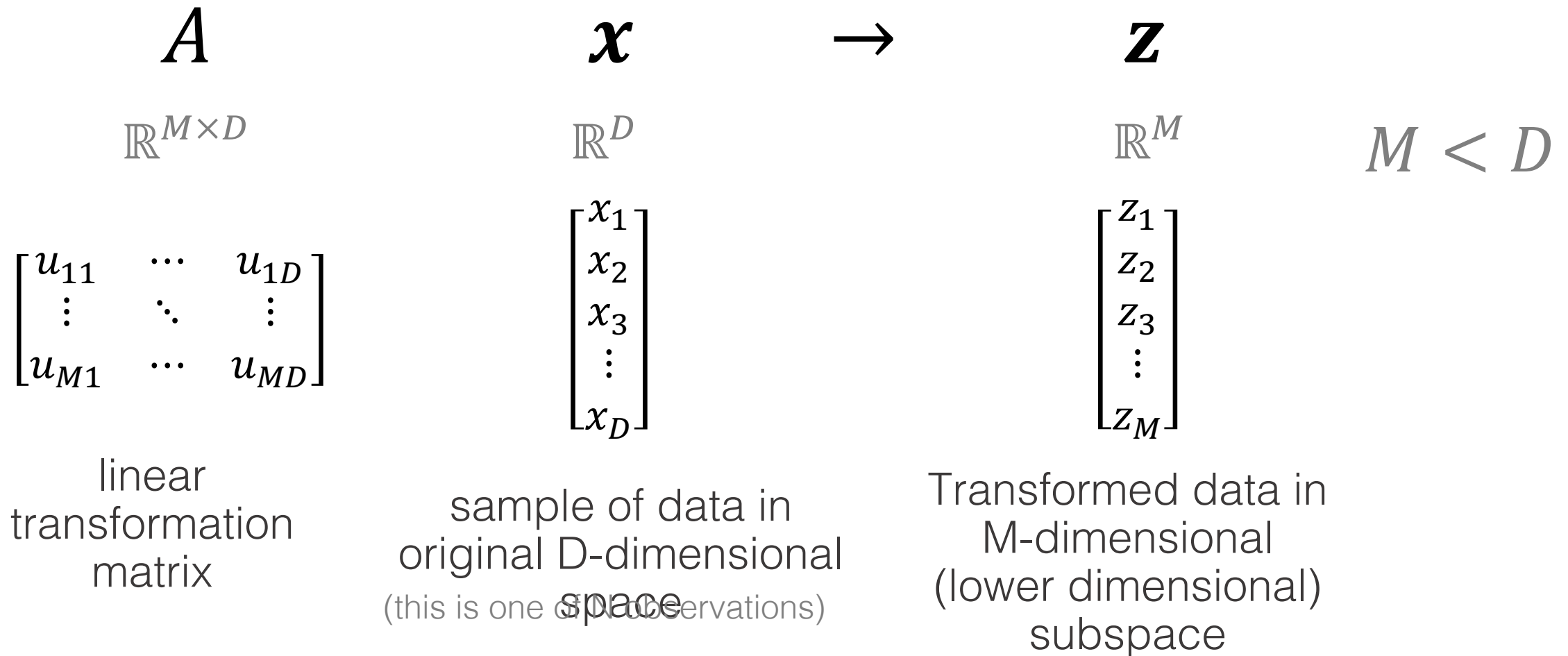
$$\mathbf{X} = \begin{bmatrix} x_{11} & \cdots & x_{1D} \\ \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{ND} \end{bmatrix} \text{ rows = observations}$$

We normalize each of the columns

Principal components analysis

a.k.a. the
Karhunen–Loève
Transform

Transform the data from a high dimensional space to a lower dimensional subspace, while minimizing the projection error



Principal components analysis

A

$$\begin{bmatrix} u_{11} & \cdots & u_{1D} \\ \vdots & \ddots & \vdots \\ u_{M1} & \cdots & u_{MD} \end{bmatrix} = \begin{bmatrix} - & \mathbf{u}_1^T & - \\ & \vdots & \\ - & \mathbf{u}_M^T & - \end{bmatrix}$$

linear
transformation
matrix

Each \mathbf{u}_i vector
represents

The i^{th} principal component:

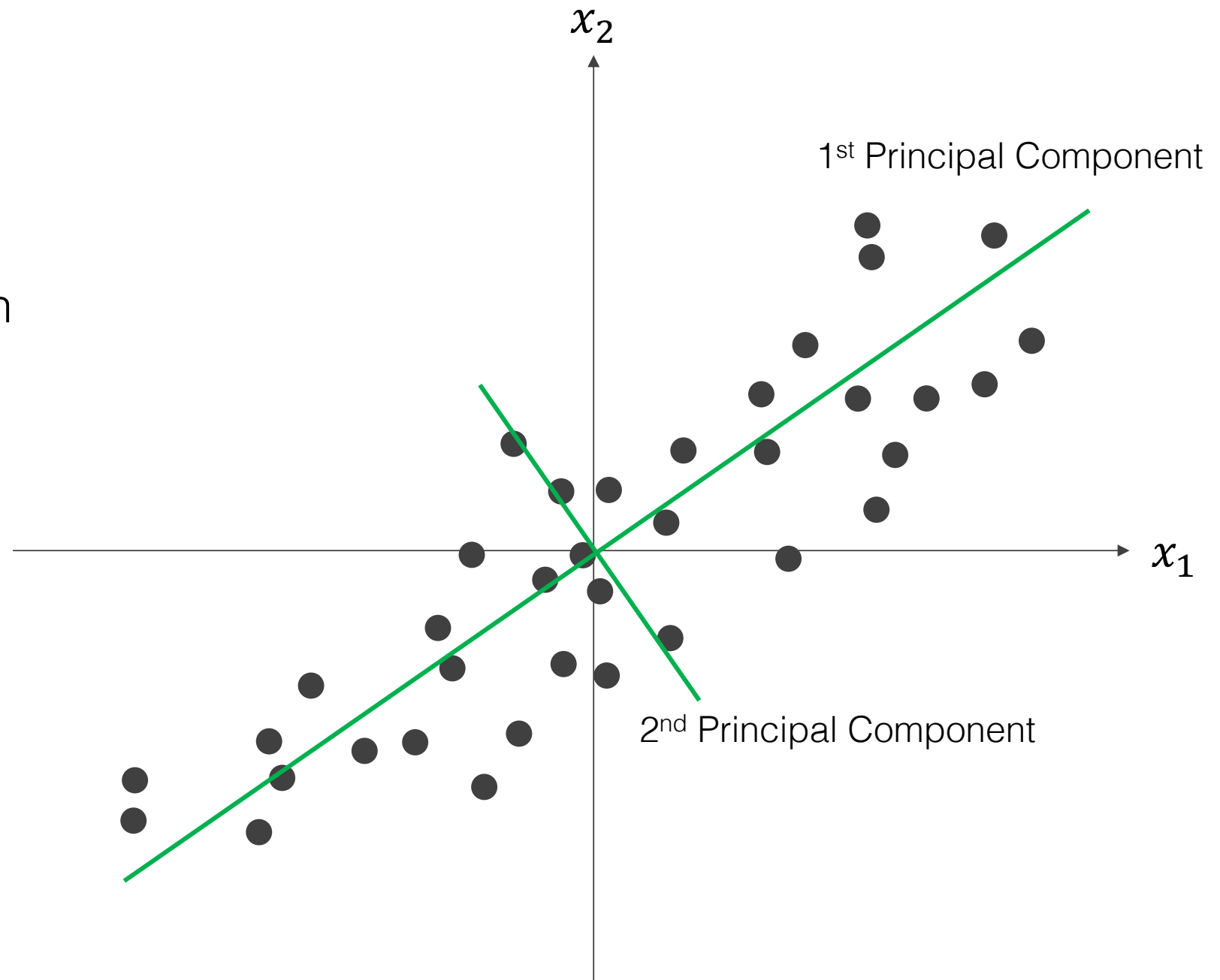
$$z_i = \mathbf{u}_i^T \mathbf{x}$$

Since only direction matters, we
assume the \mathbf{u}_i are unit vectors

$$\mathbf{u}_i^T \mathbf{u}_i = 1$$

Principal Components

Maximum variance formulation

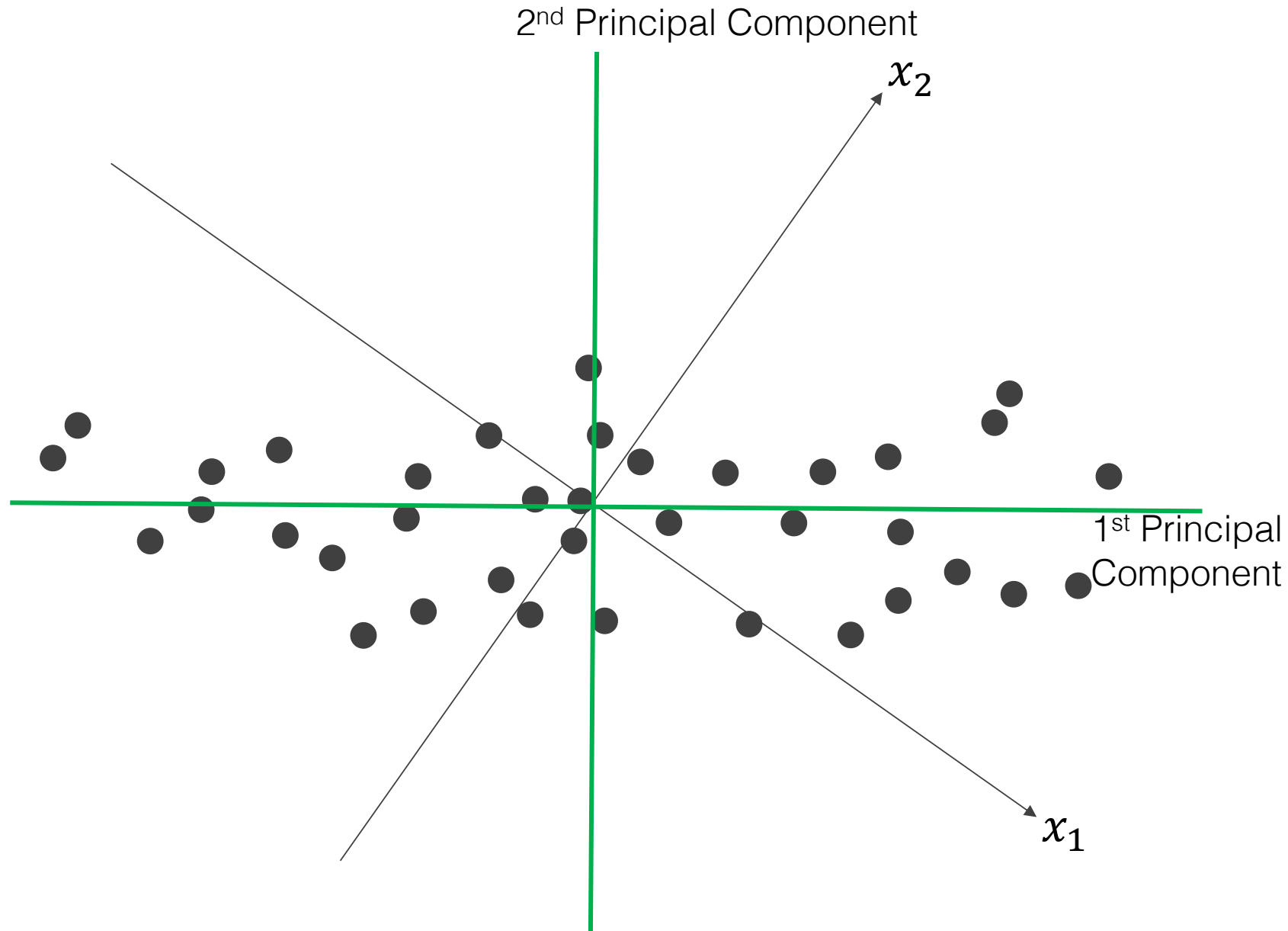


Reprojected Data onto Principal Components

Any point \mathbf{x}_j can be represented as a combination of the principle components

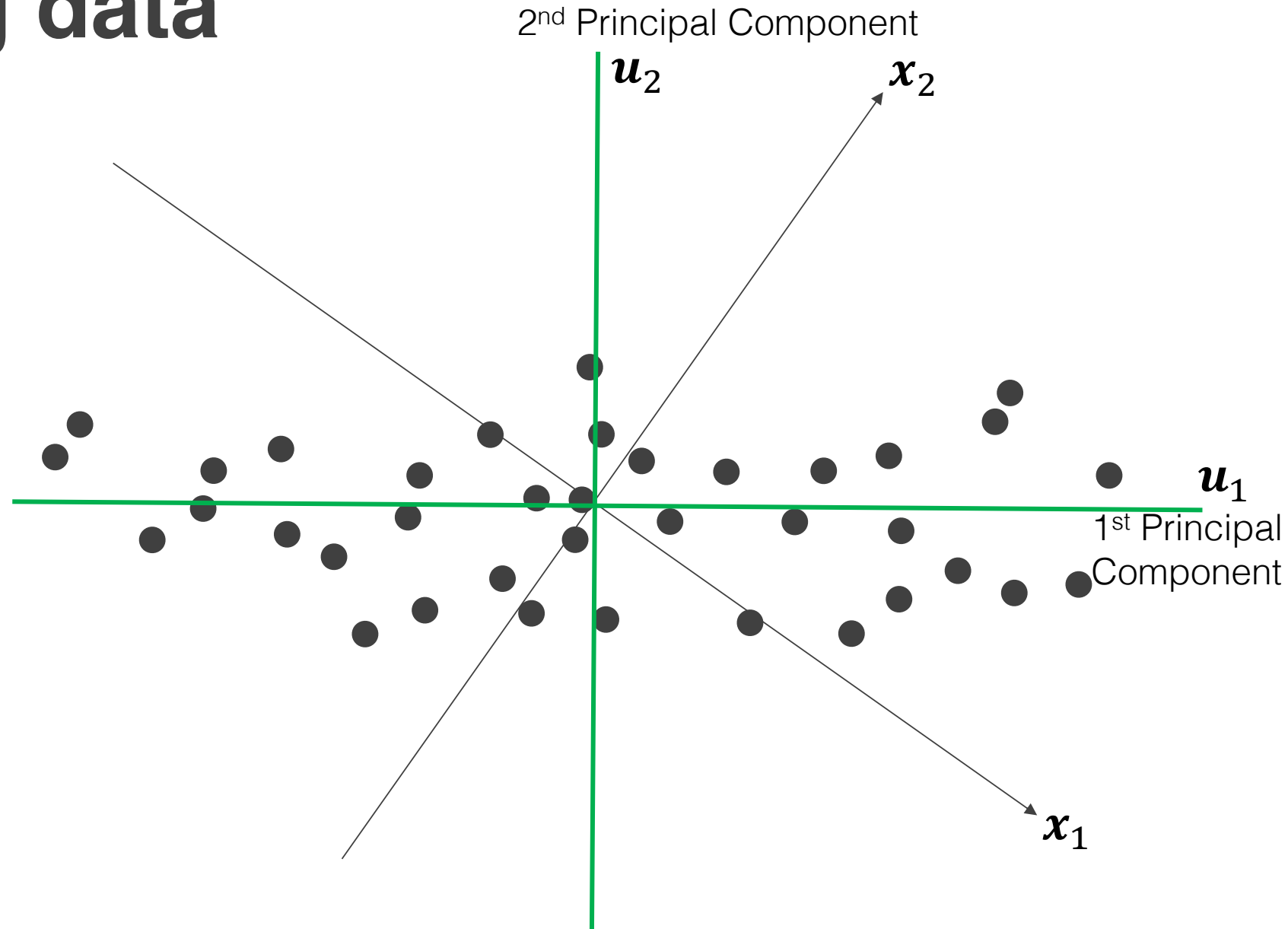
$$\mathbf{x}_j = \sum_{i=1}^D \beta_{ji} \mathbf{u}_i$$

The \mathbf{u}_i 's are an orthogonal basis for the space \mathbb{R}^D



Approximating data with principal components

$$\mathbf{x}_j = \sum_{i=1}^D (\underbrace{\mathbf{x}_j^T \mathbf{u}_i}_{\text{scalar projection}}) \underbrace{\mathbf{u}_i}_{\text{principal component direction}}$$



PCA

We want to **maximize the variance** of the projected data

Let's start by finding the **unit vector in the direction of greatest variation** in the dataset

Here the magnitude is unimportant, but the direction matters

We seek to project each point \mathbf{x}_i onto a unit PC vector. $z_i = \mathbf{u}_i^T \mathbf{x}_i$

PCA Example: find the first principal component

Mean of the data:

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \quad \mathbf{x}_i \quad [D \times 1]$$

The projected mean of the data:

$$\bar{z} = \mathbf{u}_1^T \bar{\mathbf{x}}$$

We can compute the sample variance as:

$$\sigma_z^2 = \frac{1}{N} \sum_{i=1}^N (z_i - \bar{z})^2$$

[scalar]

$$\sigma_z^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{u}_1^T \mathbf{x}_i - \mathbf{u}_1^T \bar{\mathbf{x}})^2$$

PCA Example: find the first principal component

We can compute the sample variance as:

$$\sigma_z^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{u}_1^T \mathbf{x}_i - \mathbf{u}_1^T \bar{\mathbf{x}})^2$$

$$= \frac{1}{N} \sum_{i=1}^N \mathbf{u}_1^T (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{u}_1$$

$$= \mathbf{u}_1^T \mathbf{\Sigma} \mathbf{u}_1 \quad \text{Variance of the projected data}$$

Define:

$$\mathbf{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

Covariance matrix

$$\mathbf{u}_1^T \Sigma \mathbf{u}_1 = \frac{1}{N} \sum_{i=1}^N \mathbf{u}_1^T (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{u}_1$$

$$\Sigma = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} & \cdots & \Sigma_{1D} \\ \Sigma_{21} & \Sigma_{22} & \cdots & \Sigma_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{D1} & \Sigma_{D2} & \cdots & \Sigma_{DD} \end{bmatrix}$$

$$\begin{aligned} \Sigma_{ij} &= E[(X_i - \mu_i)(X_j - \mu_j)] \\ &= \text{cov}(X_i, X_j) \end{aligned}$$

$$\text{If } \mu_i = 0 \forall i$$

$$\Sigma_{ij} = E[X_i X_j] = \frac{1}{N} \mathbf{x}_i^T \mathbf{x}_j$$

$$\Sigma = \frac{1}{N} \mathbf{X}^T \mathbf{X}$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \Sigma_{12} & \cdots & \Sigma_{1D} \\ \Sigma_{21} & \sigma_2^2 & \cdots & \Sigma_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{D1} & \Sigma_{D2} & \cdots & \Sigma_D^2 \end{bmatrix}$$

$$\sigma_i^2 = E[(X_i - \mu_i)^2]$$

Covariance matrix properties

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_1^2 & \Sigma_{12} & \cdots & \Sigma_{1D} \\ \Sigma_{21} & \sigma_2^2 & \cdots & \Sigma_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{D1} & \Sigma_{D2} & \cdots & \Sigma_D^2 \end{bmatrix}$$

Positive semidefinite and symmetric ($\mathbf{\Sigma} = \mathbf{\Sigma}^T$)

All eigenvalues are positive

Eigenvectors are orthogonal

If the features, x_1, x_2, \dots, x_D are independent, $\mathbf{\Sigma}$ is diagonal because $cov(x_i, x_j) = 0$ if $i \neq j$

PCA

We want to **maximize variance** $\sigma_z^2 = \mathbf{u}_1^T \Sigma \mathbf{u}_1$
subject to $\mathbf{u}_1^T \mathbf{u}_1 = 1$

We can use **Lagrange multipliers**:

Maximize $f(x)$

subject to the constraint $g(x)$

We maximize this: $L(x, \lambda) = f(x) - \lambda g(x)$

$$f(\mathbf{x}, \mathbf{u}_i) = \mathbf{u}_1^T \Sigma \mathbf{u}_1$$

$$g(\mathbf{x}, \mathbf{u}_i) = \mathbf{u}_1^T \mathbf{u}_1 - 1$$

For our case: $L(\mathbf{x}, \mathbf{u}_1, \lambda) = \mathbf{u}_1^T \Sigma \mathbf{u}_1 - \lambda(\mathbf{u}_1^T \mathbf{u}_1 - 1)$

We take the derivative and set it equal to zero

PCA

$$L(\mathbf{x}, \mathbf{u}_1, \lambda) = \mathbf{u}_1^T \mathbf{\Sigma} \mathbf{u}_1 - \lambda(\mathbf{u}_1^T \mathbf{u}_1 - 1)$$

We take the derivative with respect to \mathbf{u}_1 and set it equal to zero

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{u}_1} &= \frac{\partial}{\partial \mathbf{u}_1} \mathbf{u}_1^T \mathbf{\Sigma} \mathbf{u}_1 - \frac{\partial}{\partial \mathbf{u}_1} \lambda(\mathbf{u}_1^T \mathbf{u}_1 - 1) \\ &= 2\mathbf{\Sigma} \mathbf{u}_1 - 2\lambda \mathbf{u}_1 = 0 \quad (\text{since } \mathbf{\Sigma} \text{ is symmetric}) \end{aligned}$$

$\mathbf{\Sigma} \mathbf{u}_1 = \lambda \mathbf{u}_1 \quad \rightarrow \quad \mathbf{u}_1$ is an eigenvector of the covariance matrix $\mathbf{\Sigma}$, and λ is an eigenvalue

PCA

Since we want to maximize the variance in the projected features:

We want to maximize: $\sigma_z^2 = \mathbf{u}_1^T \mathbf{\Sigma} \mathbf{u}_1$

And we know that: $\mathbf{\Sigma} \mathbf{u}_1 = \lambda \mathbf{u}_1$

So we can write: $\sigma_z^2 = \mathbf{u}_1^T \lambda \mathbf{u}_1 = \lambda \mathbf{u}_1^T \mathbf{u}_1 = \lambda$

Therefore we choose the eigenvector that corresponds to the **largest eigenvalue**

PCA: Variance explained

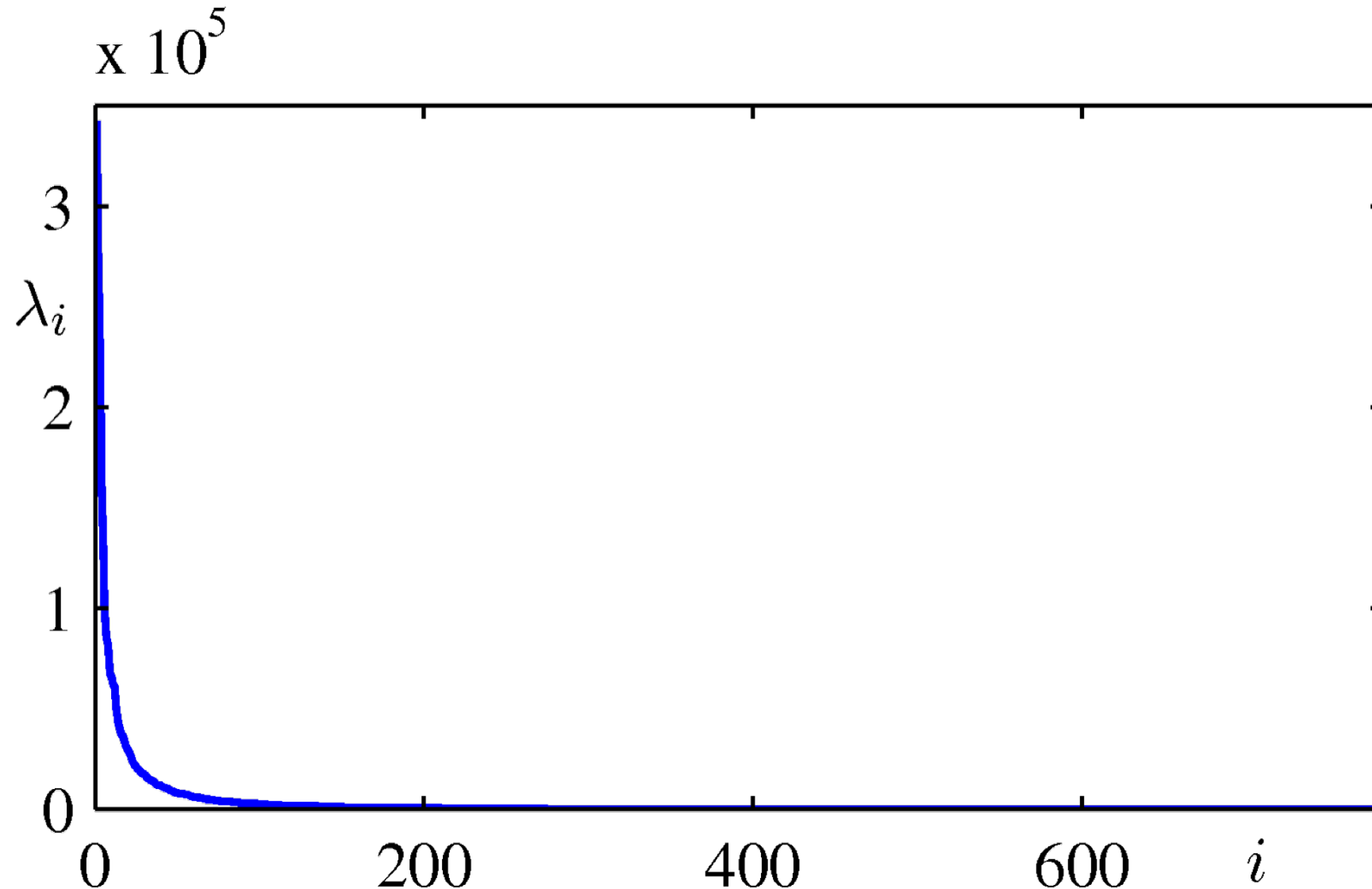
The **fraction of variance explained** =
$$\frac{\sum_{i=1}^M \lambda_i}{\sum_{i=1}^D \lambda_i}$$

M = dimensionality of the subspace

D = dimensionality of the original data space

The more principle components included, the more of the variance will be represented in the projected data

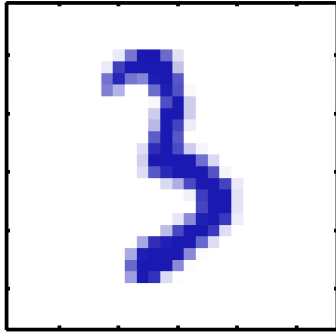
Eigenvalues by principal component i



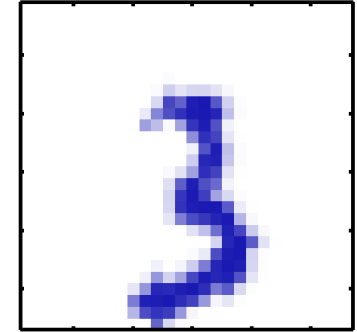
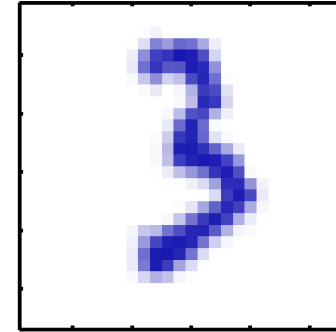
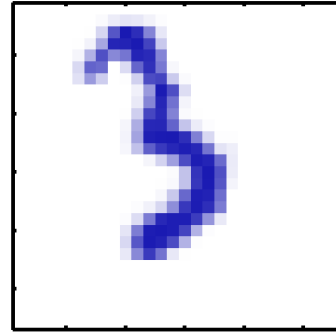
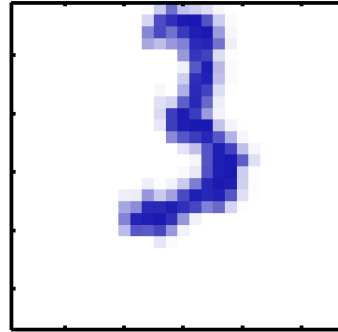
Bishop, Pattern Recognition, 2006

Example: translated digits

Original digit



Translated digits



Types of translation:

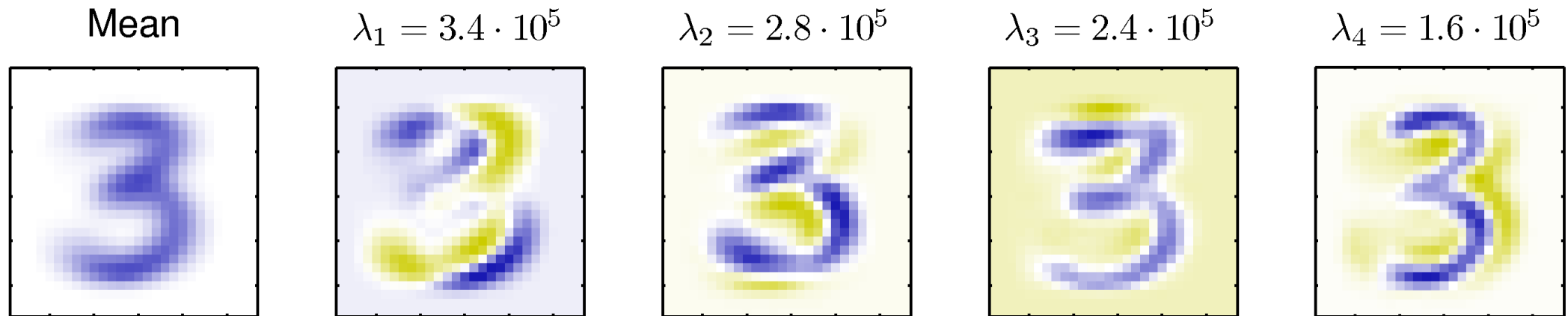
1. Horizontal translation
2. Vertical translation
3. Rotation

Original digits: 64 x 64 pixels

New size: 100 x 100 pixels

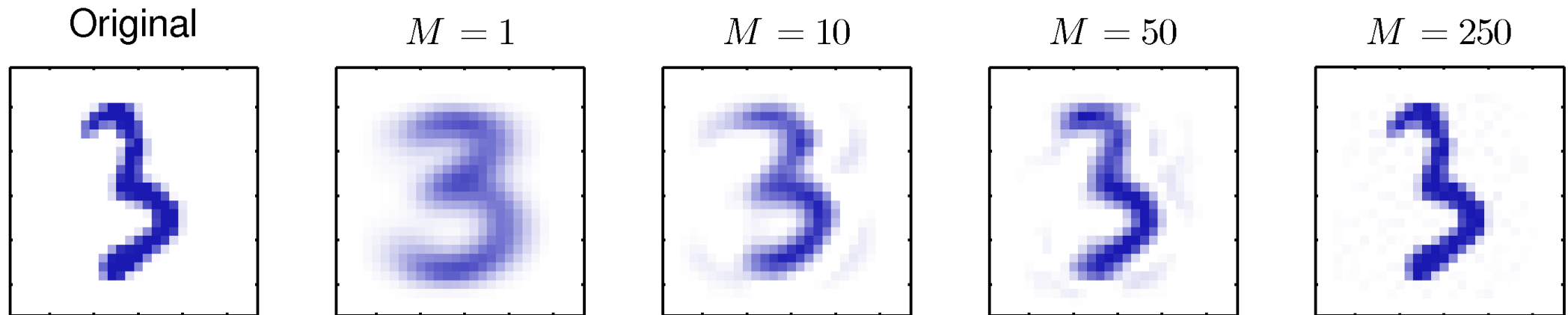
Example: translated digits

Examples of first four principle component eigenvectors and eigenvalues:



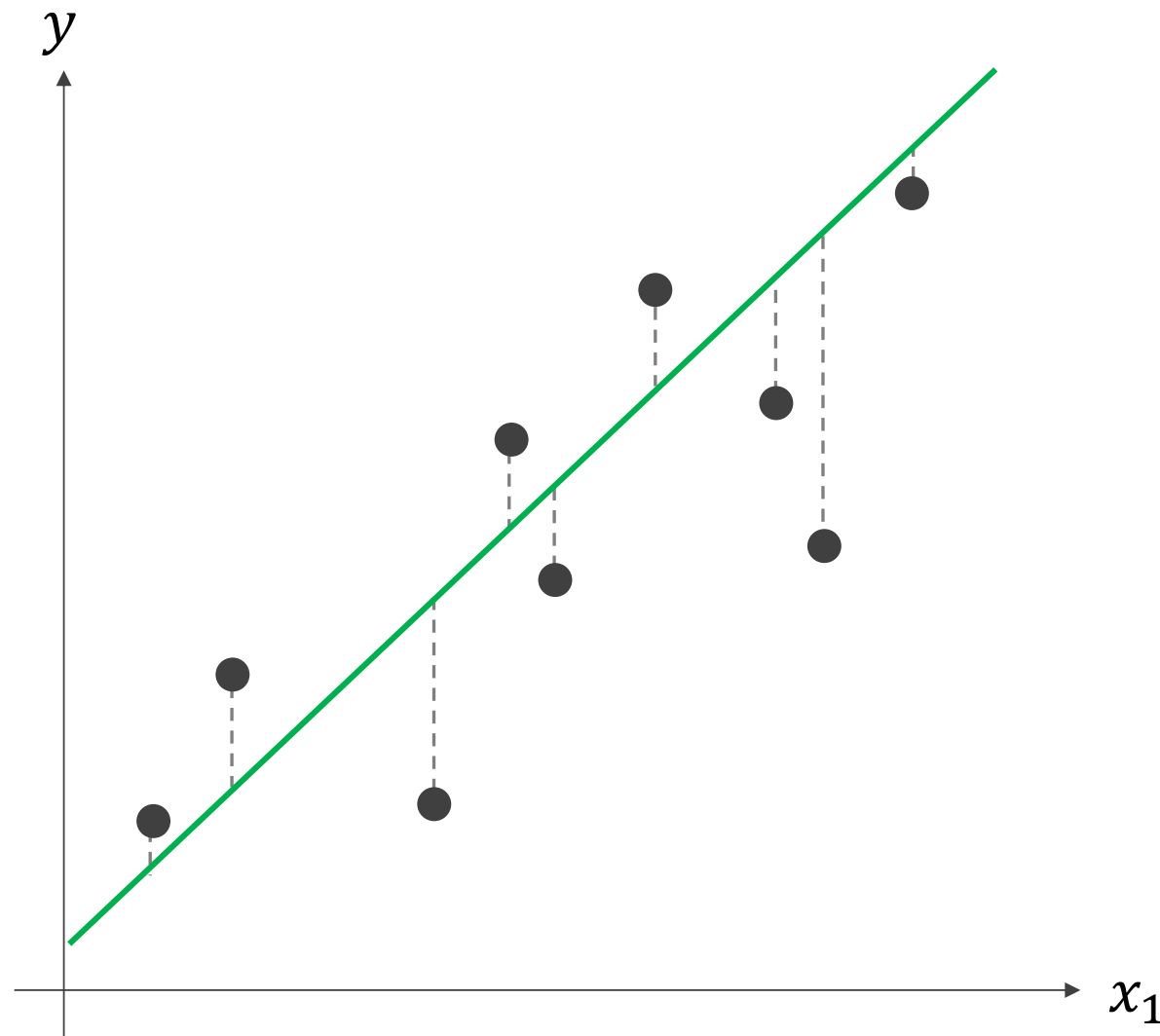
Example: translated digits

Reconstructed examples using different numbers of principal components:

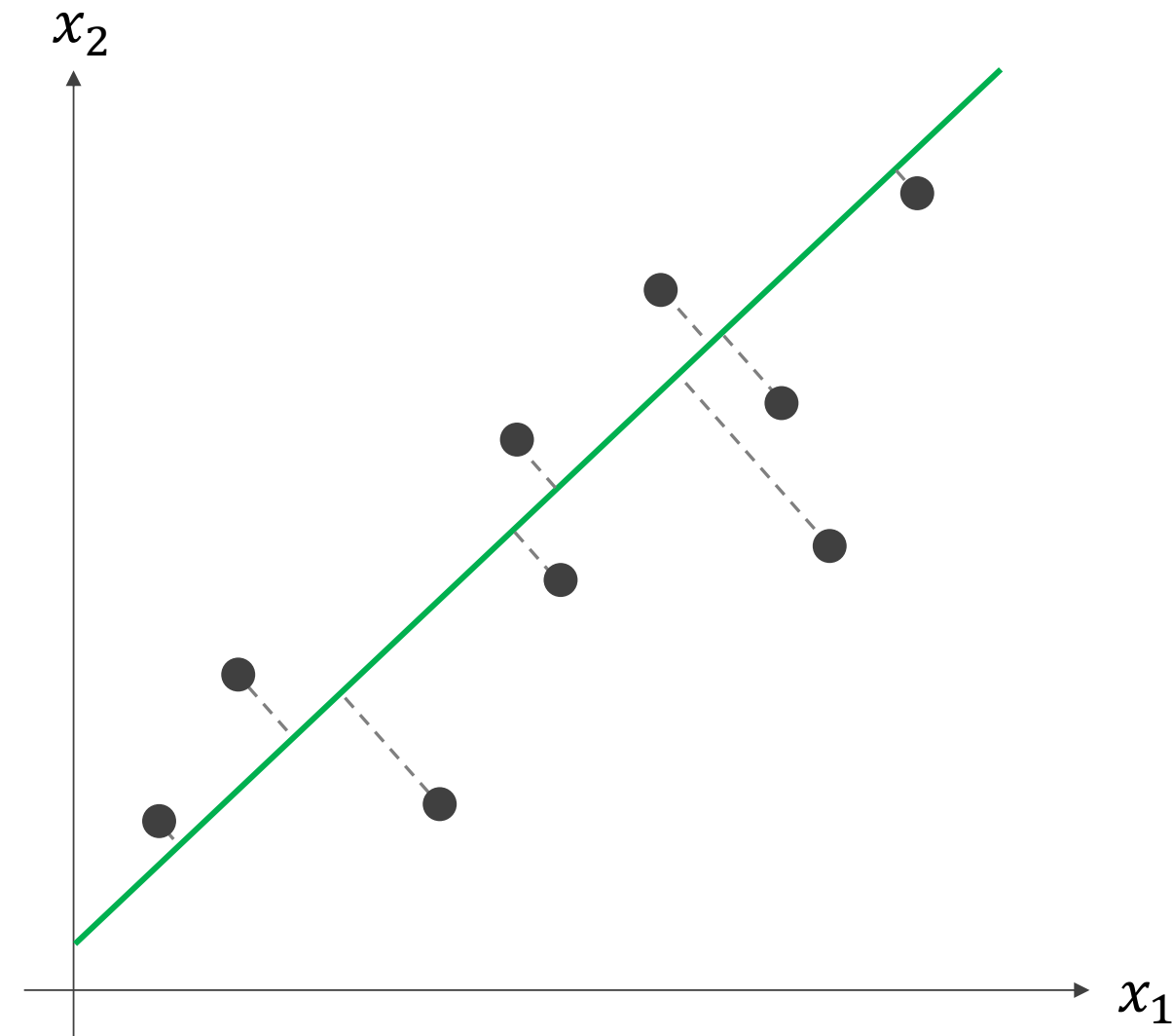


Relationship between the objective for least squares and PCA

Linear regression



First principal component



Extracting principal components

- 0 **Goal:** reduce the dimensionality of our data from D to M , where $M < D$
- 1 Normalize each feature to mean zero and a standard deviation of 1
- 2 Determine the principal components
 - Calculate the eigenvectors and eigenvalues of the data covariance matrix, Σ
 - Eigenvectors in descending order of their eigenvalues are the principal components
- 3 Project the data features on the principal components
- 4 Keep the top M principal components to reduce into a lower dimension

columns = features (D)

$$X = \begin{bmatrix} x_{11} & \cdots & x_{1D} \\ \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{ND} \end{bmatrix}$$

rows = observations (N)

size $[N \times D]$

Each observation as a vector:

$$x_i \quad i = 1, \dots, D$$

size $[D \times 1]$

u_i eigenvectors / principal components

size $[D \times 1]$

λ_i eigenvalues (how much of the variance is explained)

size [scalar]

$i = 1, \dots, D$

$$z_{ij} = u_j^T x_i \quad \begin{matrix} j = 1, \dots, D \\ i = 1, \dots, N \end{matrix}$$

size [scalar]

$$A = [u_1, u_2, \dots, u_M]$$

size $[D \times M]$

$$z_i = A^T x_i \quad i = 1, \dots, N$$

example

Each pixel represents a feature

$x_1 \quad x_2 \quad x_3 \quad x_4$

$u_1 \quad u_2 \quad u_3 \quad u_4$

$u_1 \cdot x_1 = z_{11}$

$[D \times 1] \cdot [D \times 1] = 2.43$ [scalar]

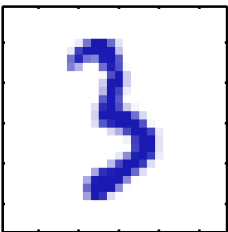
Images from Bishop, Pattern Recognition, 2006

Reconstructing our data from principal components

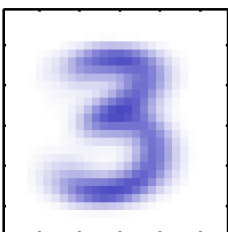
Sum the product of our projected data, \mathbf{z}_i , and our principle components

$$\hat{\mathbf{x}}_i = \sum_{j=1}^M z_{ij} \mathbf{u}_j$$

Example: the i^{th} observation: $\mathbf{x}_i =$



$\bar{\mathbf{x}} =$



PCA-projected data: $\mathbf{z}_i = [z_{i1}, z_{i2}, \dots, z_{iM}]$

$M = 1$

$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}} + z_{i1} \cdot \mathbf{u}_1 =$$

The diagram shows the reconstruction of the digit '3' using only the mean image $\bar{\mathbf{x}}$ and the first principal component \mathbf{u}_1 . The result is a blurred version of the digit.

$M = 250$

$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}} + z_{i1} \cdot \mathbf{u}_1 + z_{i2} \cdot \mathbf{u}_2 + \sum_{j=3}^{250} z_{ij} \mathbf{u}_j =$$

The diagram shows the reconstruction of the digit '3' using the mean image $\bar{\mathbf{x}}$ and the first 250 principal components $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{250}$. The result is a sharper version of the digit.

$M = 10,000$

$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}} + z_{i1} \cdot \mathbf{u}_1 + z_{i2} \cdot \mathbf{u}_2 + \sum_{j=3}^{10,000} z_{ij} \mathbf{u}_j =$$

The diagram shows the reconstruction of the digit '3' using the mean image $\bar{\mathbf{x}}$ and the first 10,000 principal components $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{10,000}$. The result is a very sharp, clear version of the digit.

(perfect reconstruction)

Why PCA?

- Dimensionality reduction
- Feature extraction
- Data visualization
- Lossy data compression

Other dimensionality reduction techniques

- Kernel PCA
- Random projections
- Multidimensional scaling
- Locality sensitive hashing
- Autoencoders
- Isomap

e.g. Manifold Learning

