

Dimensionality Reduction

Lecture 12

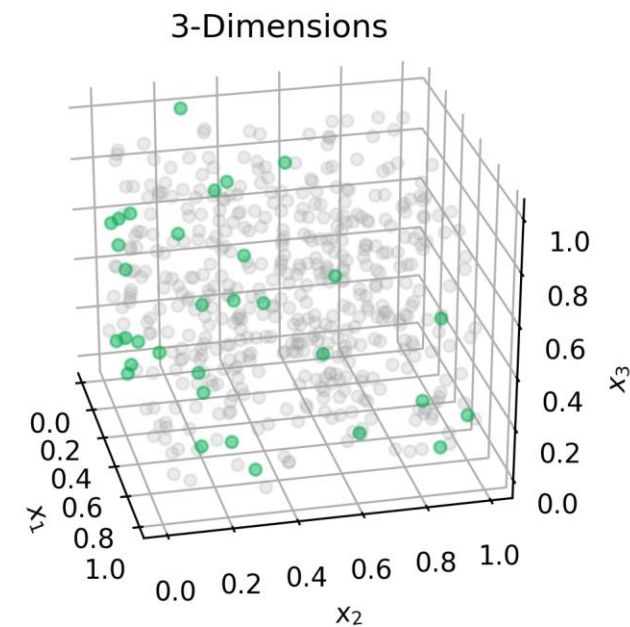
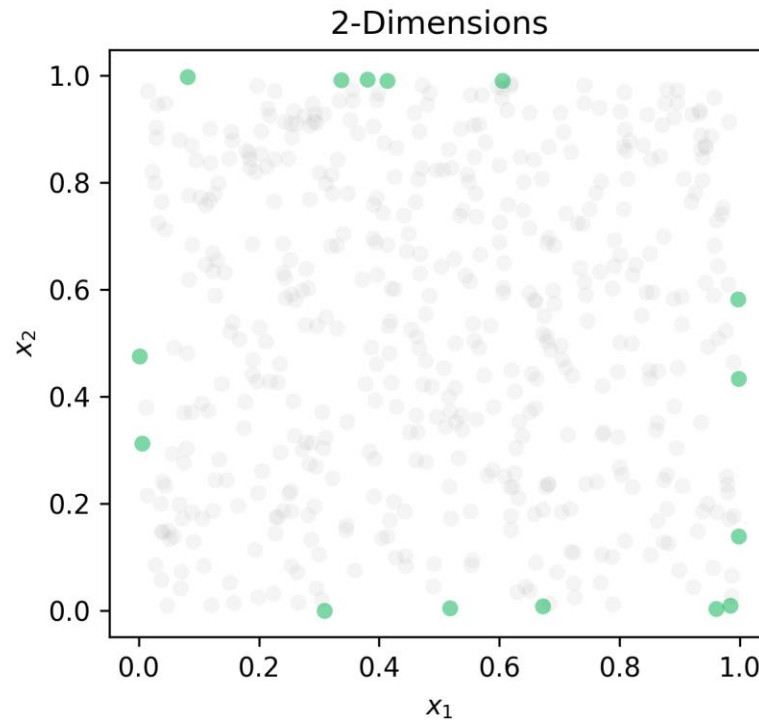
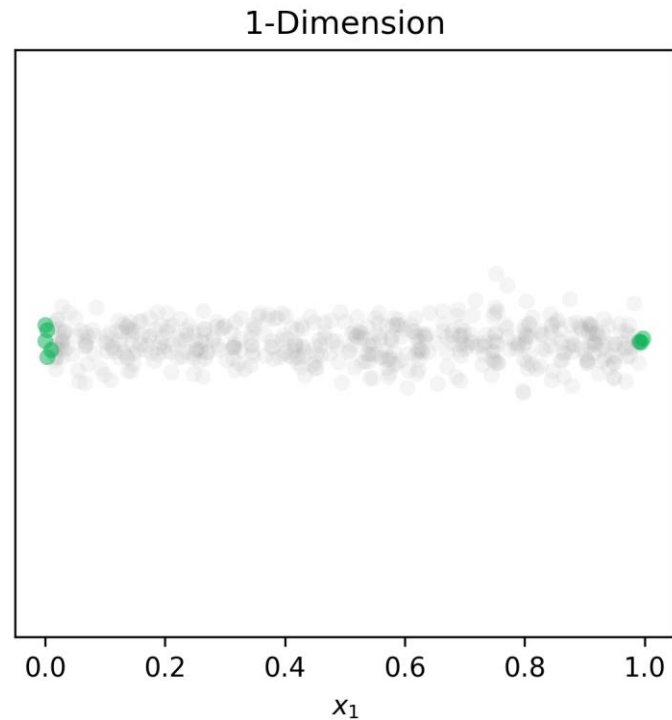
The **Curse** of Dimensionality

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

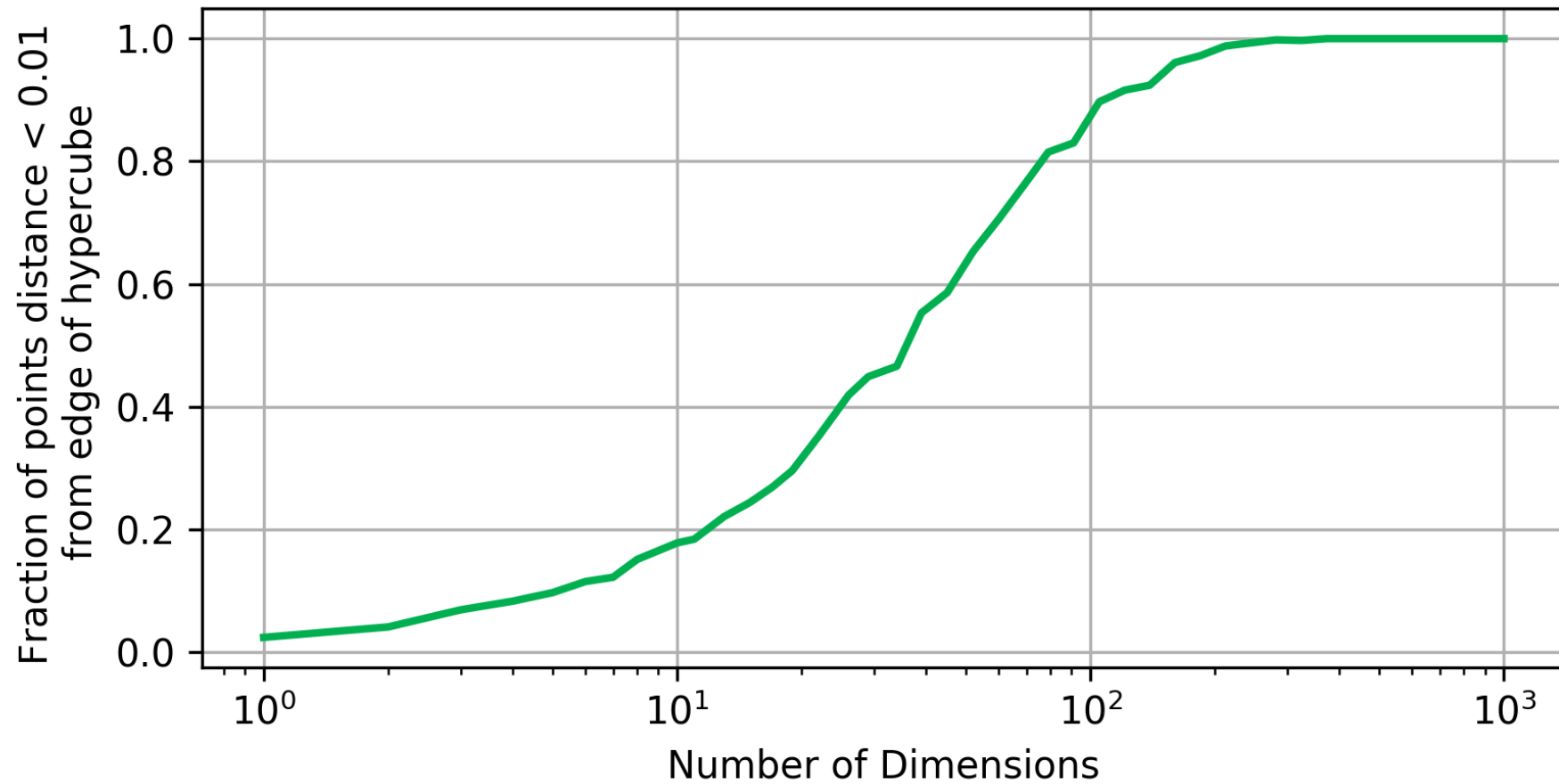
● + ● =

0.016

0.030

0.064

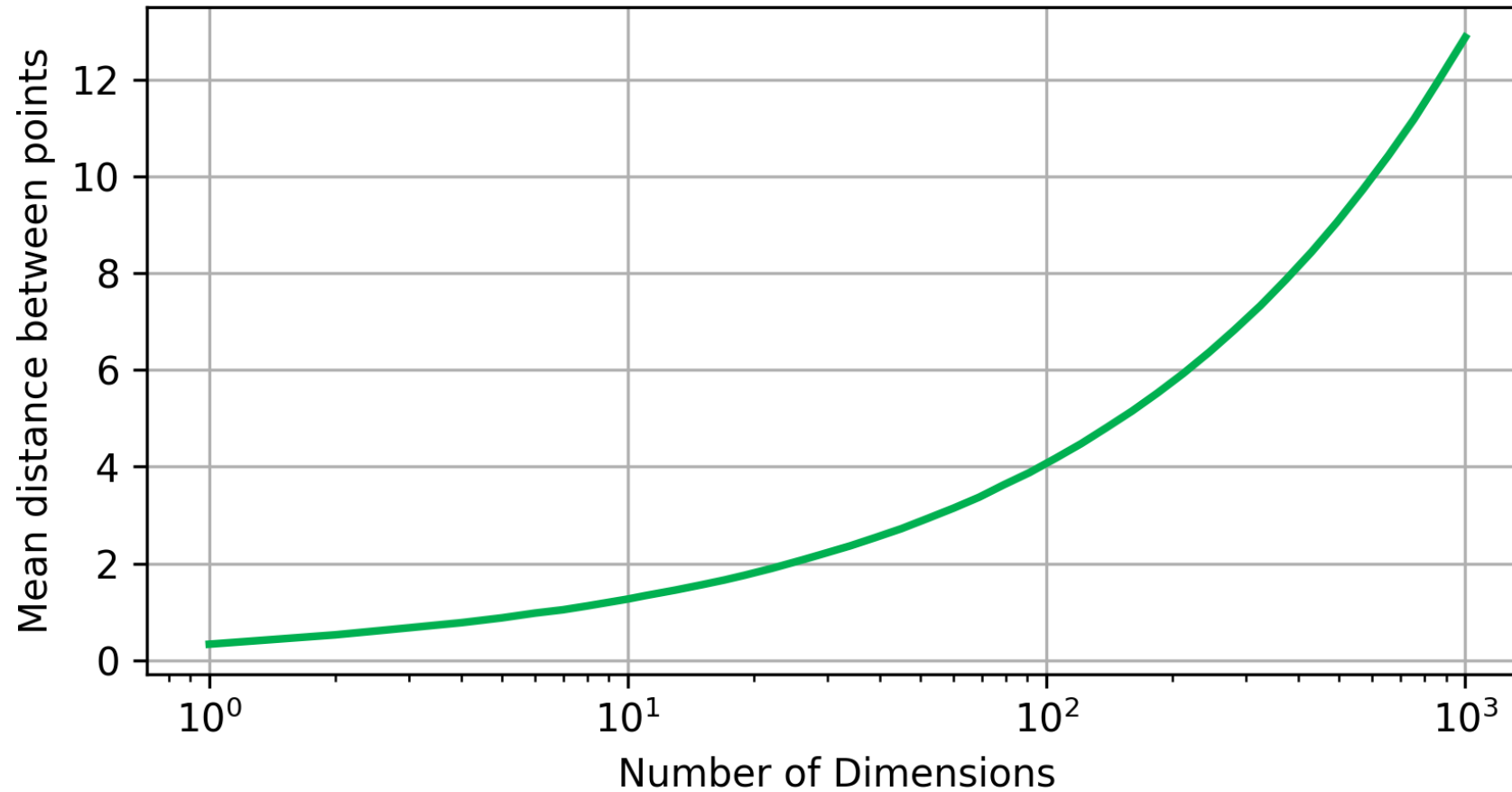
In high dimensions...



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

Note: figures constructed using 1,000 random points

In high dimensions...



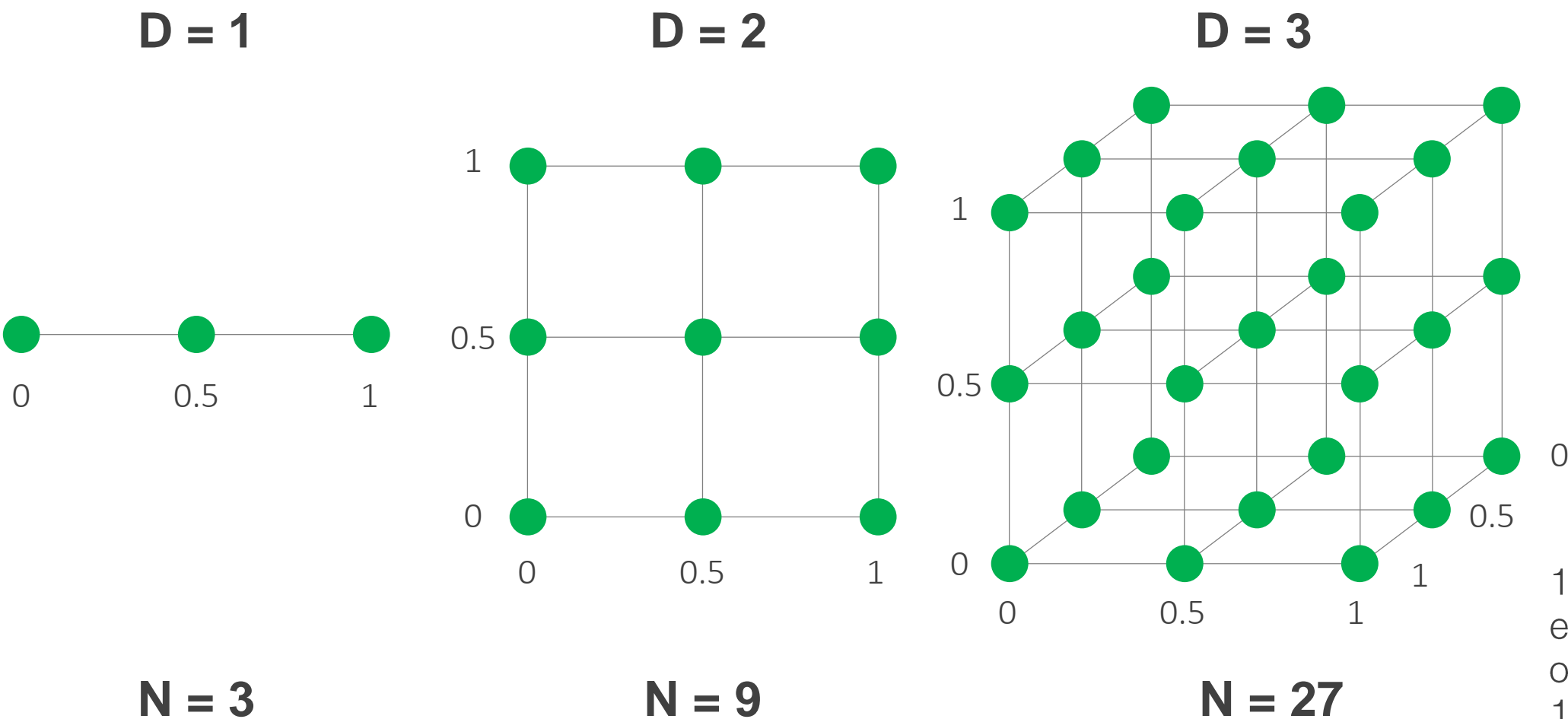
...data become sparse

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



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 computation
- Reduced redundancy of features
- Improved numerical stability due to removed correlations

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

↓ ↓

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

columns = features

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

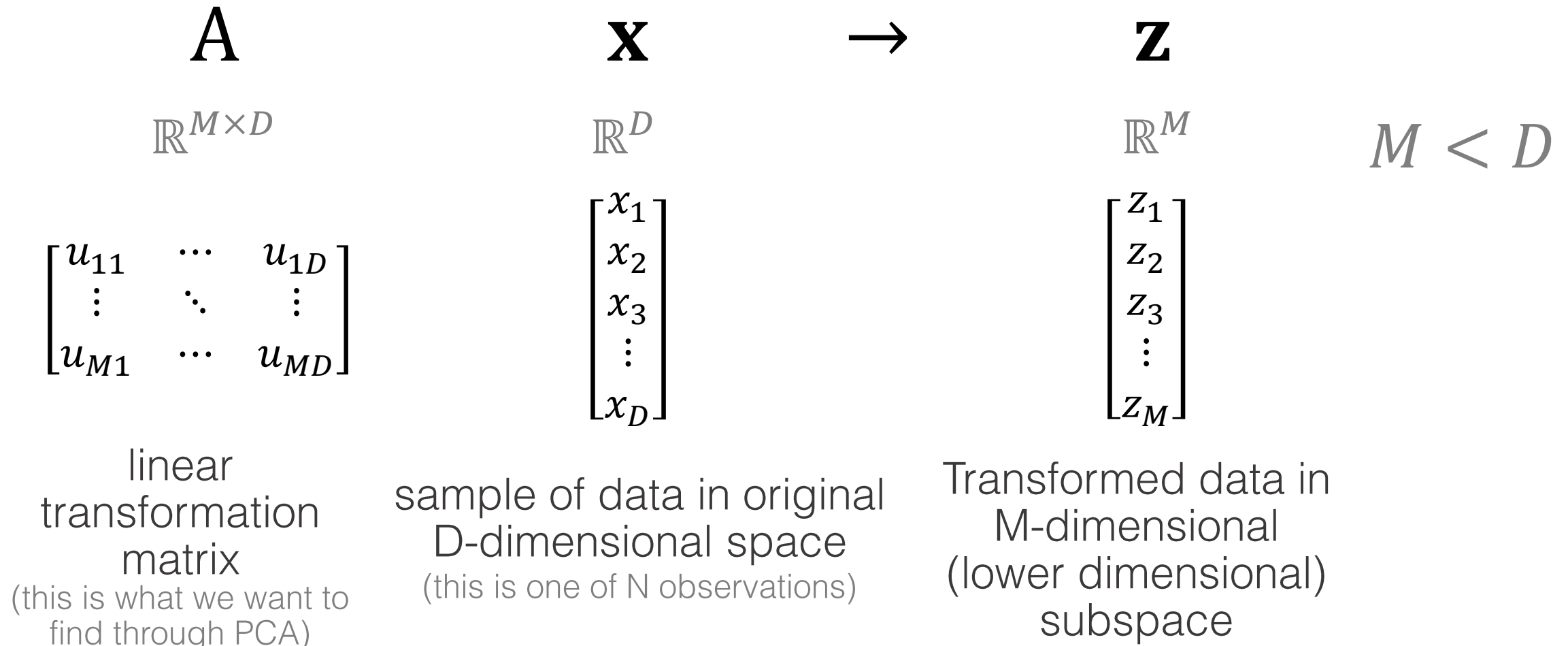
rows = observations

We normalize each of the columns

Principal components analysis

a.k.a.
Karhunen–Loève Transform
Proper orthogonal decomposition
Hotelling 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
represents a
unit vector

The i^{th} principal component:

$$\underset{\substack{\uparrow \\ \text{scalar}}}{z_i} = \underset{\substack{\uparrow \\ \text{unit vector}}}{\mathbf{u}_i^T} \mathbf{x}$$

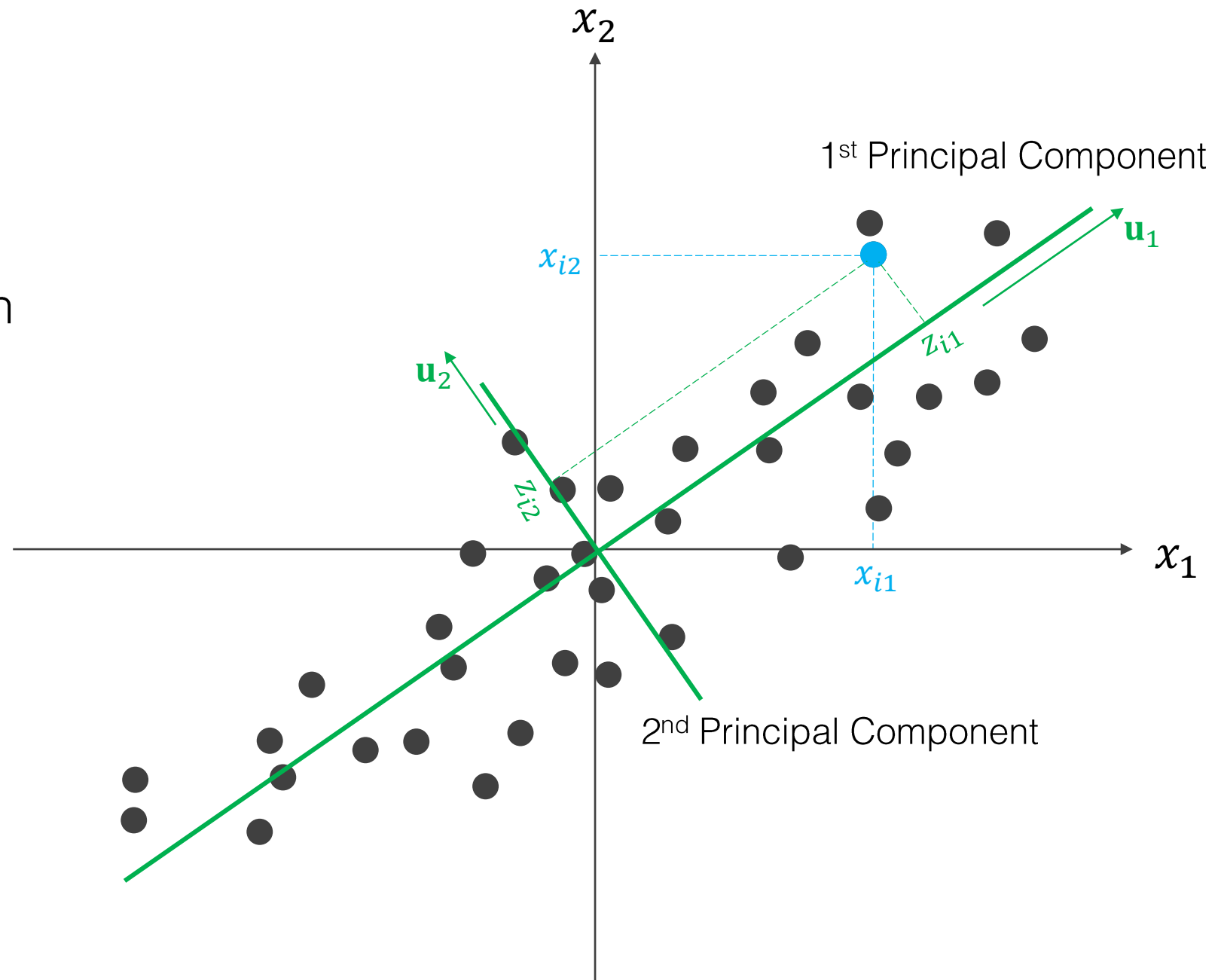
Since only direction matters, we
assume the \mathbf{u}_i are unit vectors

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

Principal Components

Maximum variance formulation

$$\mathbf{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix}$$

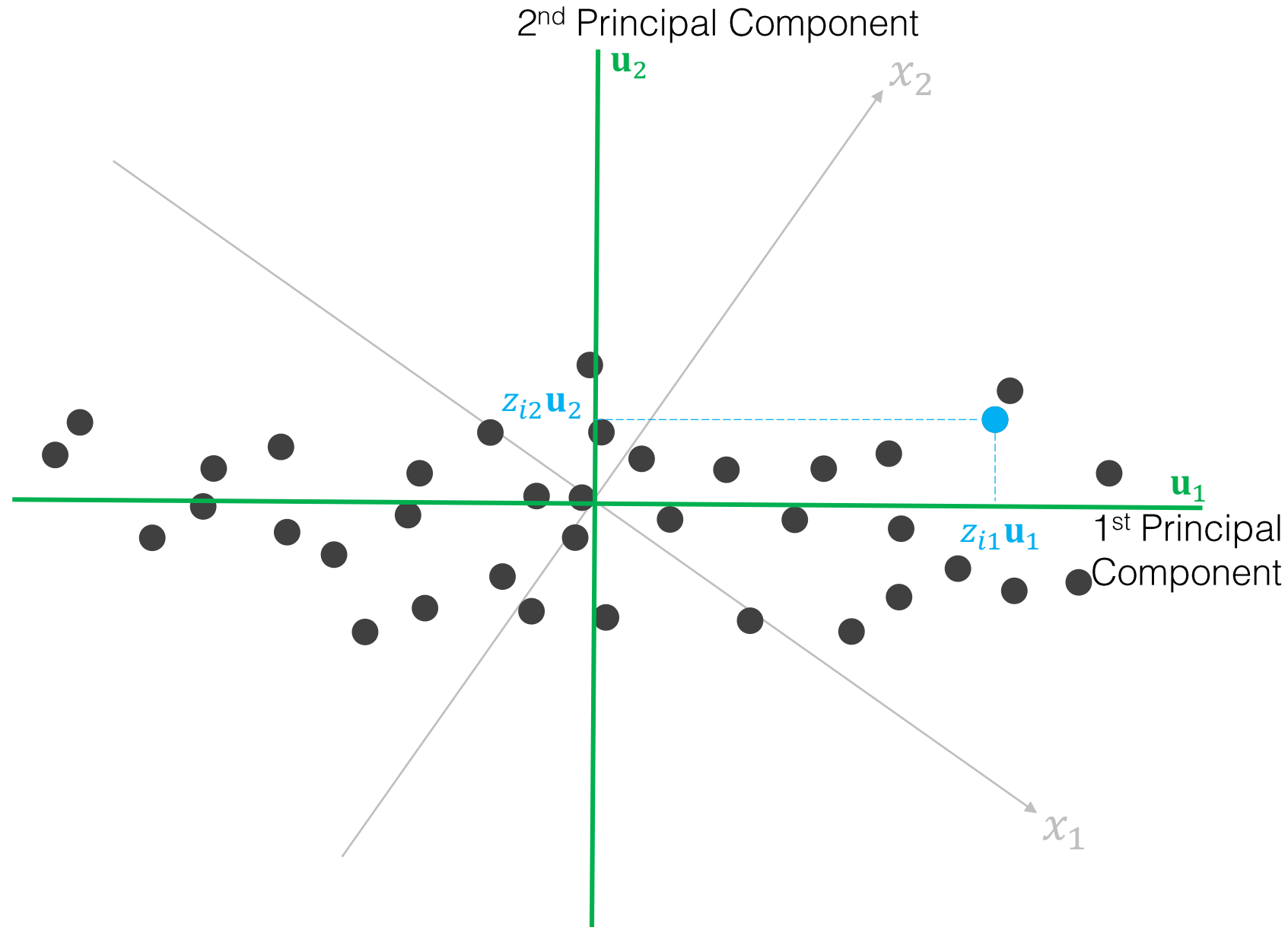


Reprojected Data onto Principal Components

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

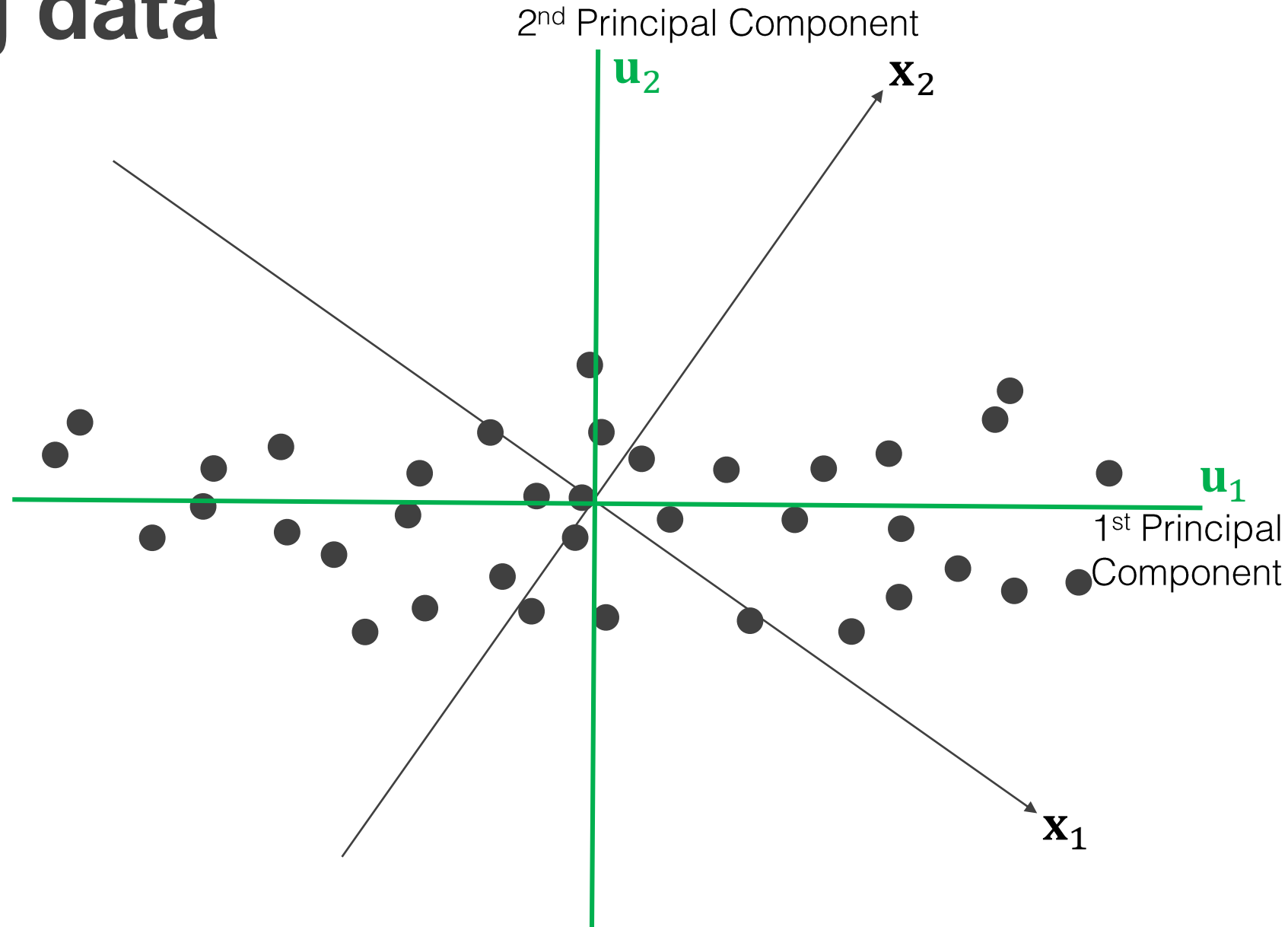
$$\mathbf{x}_i = \sum_{j=1}^D \mathbf{z}_{ij} \mathbf{u}_j$$

The \mathbf{u}_j '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

Goal

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. $z_i = \mathbf{u}_1^T \mathbf{x}_i$

PCA: Compute the variance of the transformed data

Mean of the data:

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

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]

The magnitude z_i of our data \mathbf{x}_i projected onto the unit vector \mathbf{u}_1 is:

$$z_i = \mathbf{u}_1^T \mathbf{x}_i$$

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

PCA: Compute the variance of the transformed data

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 of our data

Covariance matrix

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

$[D \times 1][1 \times D]$

$$\mathbf{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{iD} \end{bmatrix}$$

Vector of observation i

$$x_{ij}$$

Observation index Predictor index

$$\mathbf{X} = \begin{matrix} \text{Predictors} \rightarrow \\ \begin{matrix} \text{Observations} \downarrow \\ \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1D} \\ x_{21} & x_{22} & \cdots & x_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{ND} \end{bmatrix} \end{matrix} \end{matrix}$$

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

$$\begin{aligned} \Sigma_{jk} &= \frac{1}{N} \sum_{i=1}^N (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k) \\ &= \text{cov}(X_j, X_k) \\ &= E[(X_j - \mu_j)(X_k - \mu_k)] \end{aligned}$$

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

$$\begin{aligned} \sigma_j^2 &= \frac{1}{N} \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2 \\ &= E[(X_j - \mu_j)^2] \end{aligned}$$

Mean of each predictor

If $\mu_j = 0$ for all j
This will be the case IF the data are standardized

$$\begin{aligned} \Sigma_{jk} &= \frac{1}{N} \sum_{i=1}^N x_{ij} x_{ik} \\ &= \frac{1}{N} \mathbf{x}_j^T \mathbf{x}_k \end{aligned}$$

$$= E[X_j X_k]$$

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

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 ($\mathbf{v}^T \mathbf{\Sigma} \mathbf{v} \geq 0$ for all \mathbf{v}) and symmetric ($\mathbf{\Sigma} = \mathbf{\Sigma}^T$)

All eigenvalues are positive

Eigenvectors are orthogonal

If the features (predictors), x_1, x_2, \dots, x_D are independent, $\mathbf{\Sigma}$ is diagonal because $\text{cov}(X_j, X_k) = 0$ if $j \neq k$

PCA: Maximize the variance

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

We can use **Lagrange multipliers**:

Maximize $f(x)$

subject to the constraint $g(x)$

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

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

$$g(\mathbf{x}, \mathbf{u}_i) = \mathbf{u}_1^T \mathbf{u}_1 - 1 = 0$$

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

How do we know which eigenvector to use as the first principal component?

Eigenanalysis and PCA

Eigenvector Demo:

<http://setosa.io/ev/eigenvectors-and-eigenvalues/>

PCA Demo:

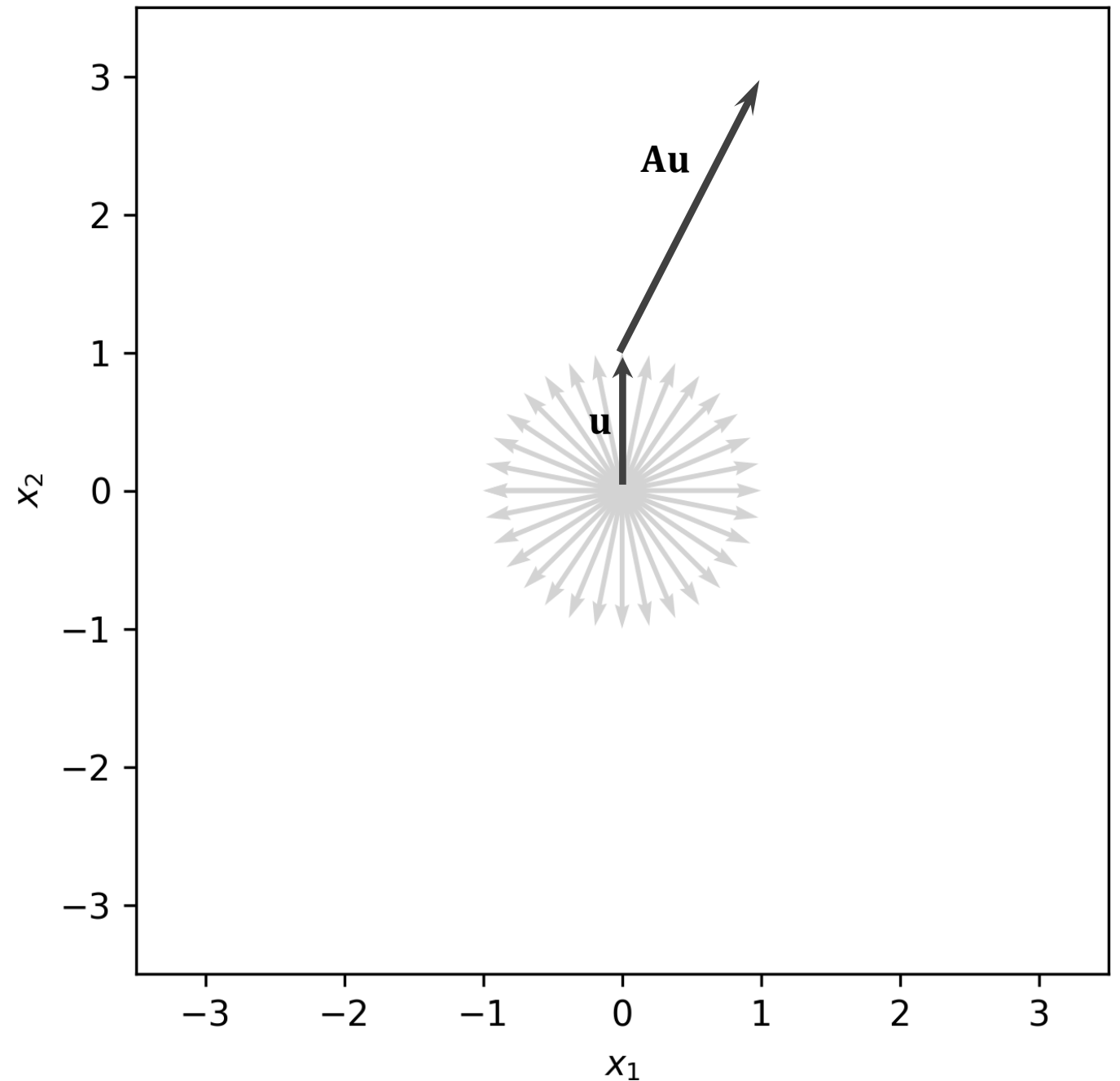
<http://setosa.io/ev/principal-component-analysis/>

Eigenvalues/Eigenvectors

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

$$\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

$$\mathbf{v} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \mathbf{A}\mathbf{v} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$



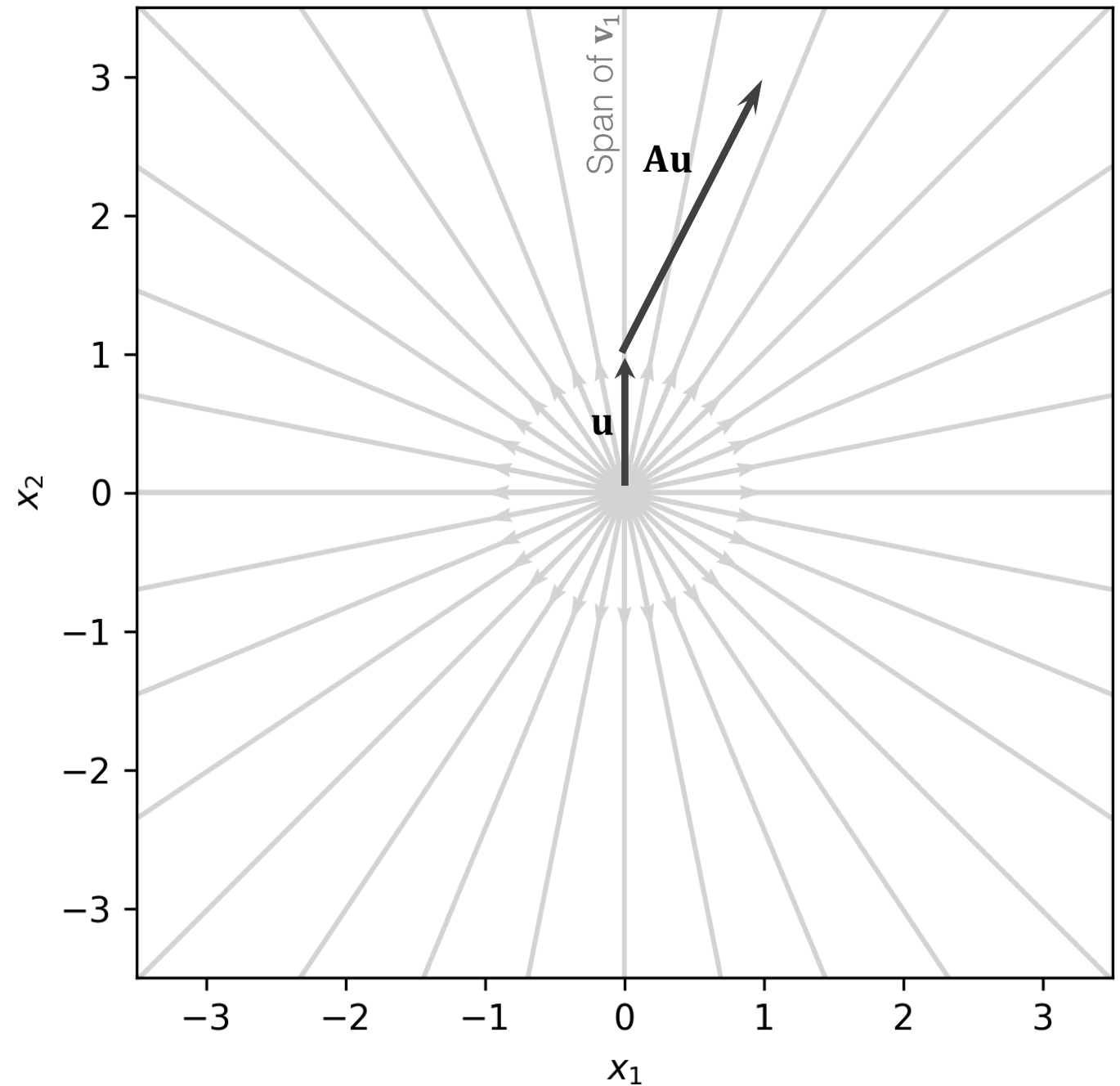
Eigenvalues/Eigenvectors

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

$$\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

$$\mathbf{u} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

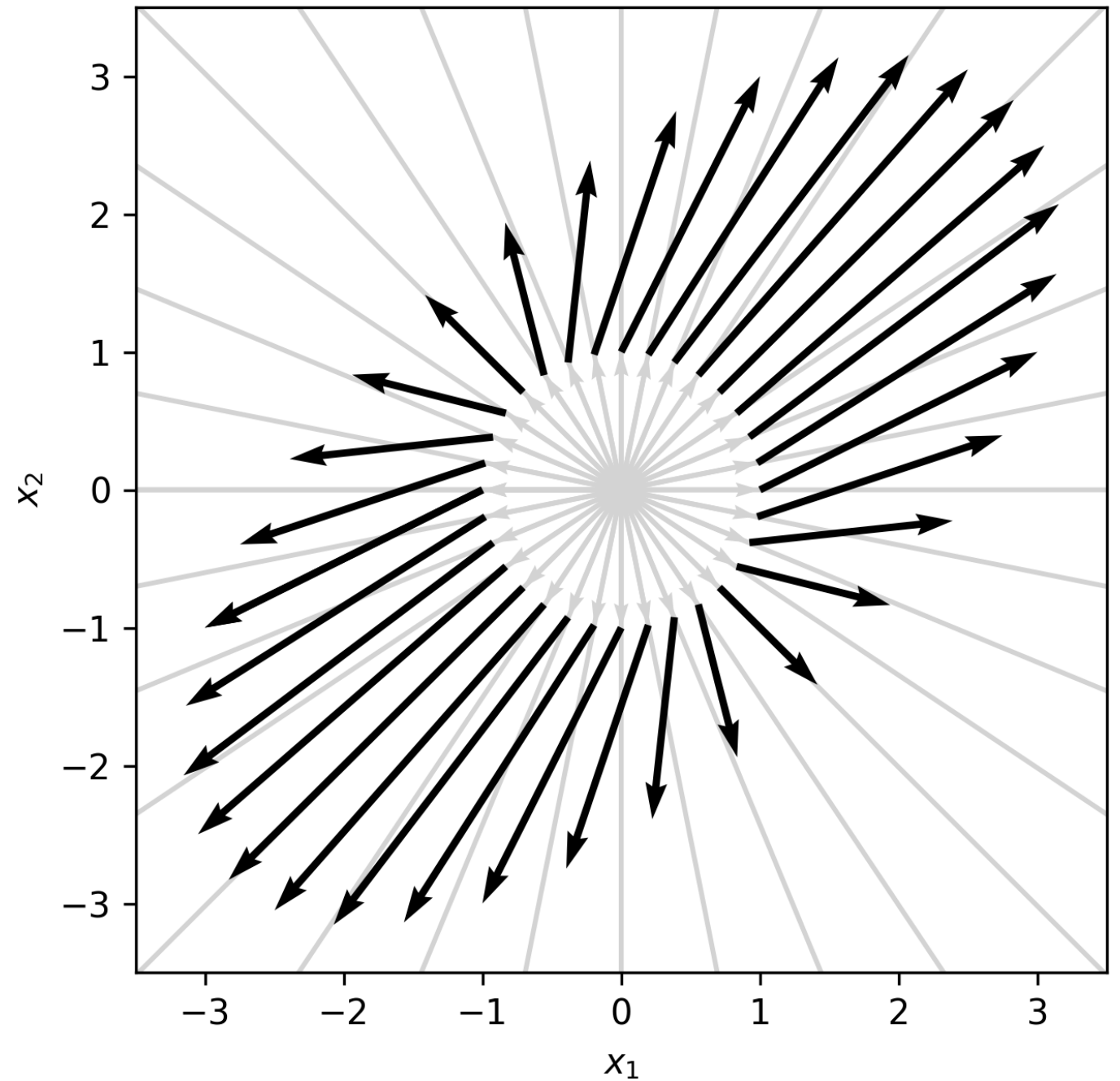
$$\begin{aligned} \mathbf{A}\mathbf{u} &= \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 2 \\ 1 \end{bmatrix} \end{aligned}$$



Eigenvalues/Eigenvectors

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

$$\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$



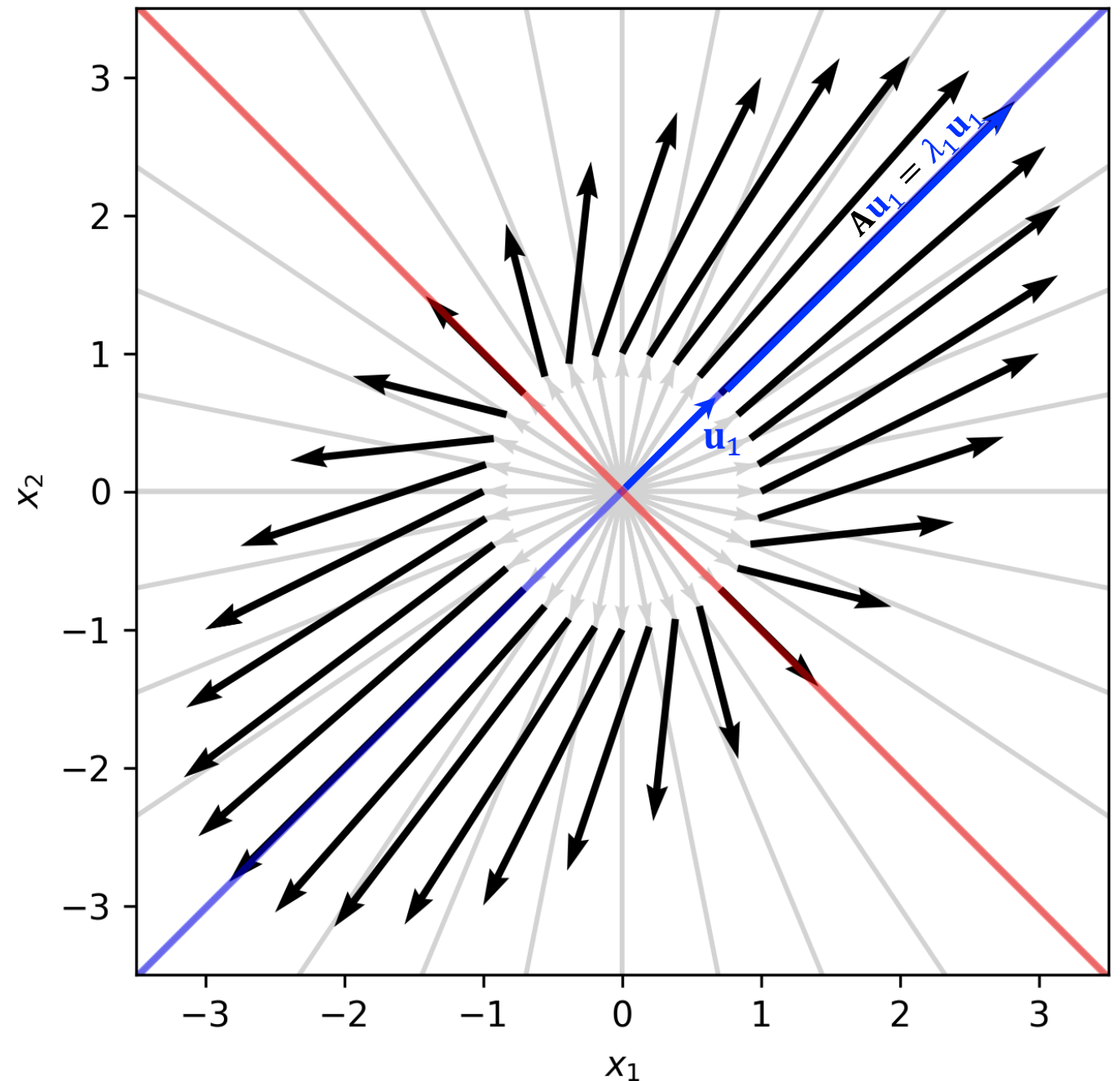
Eigenvalues/Eigenvectors

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

$$\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

$$\mathbf{u}_1 = \begin{bmatrix} 0.707 \\ 0.707 \end{bmatrix} \quad \lambda_1 = 3$$

$$\begin{aligned} \mathbf{A}\mathbf{u}_1 &= \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 0.707 \\ 0.707 \end{bmatrix} \\ &= \begin{bmatrix} 2.12 \\ 2.12 \end{bmatrix} \\ &= \lambda_1 \mathbf{u}_1 = 3 \begin{bmatrix} 0.707 \\ 0.707 \end{bmatrix} \end{aligned}$$



PCA

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

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

To do that this must be true:
(shown on last slide) $\mathbf{\Sigma} \mathbf{u}_1 = \lambda \mathbf{u}_1$

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

Variance corresponding
to our first principle
component

Therefore we choose as our first principle component the eigenvector that corresponds to the **largest eigenvalue**

The first PC will account for the most variance, the second PC to the second most, etc.

PCA: Variance explained

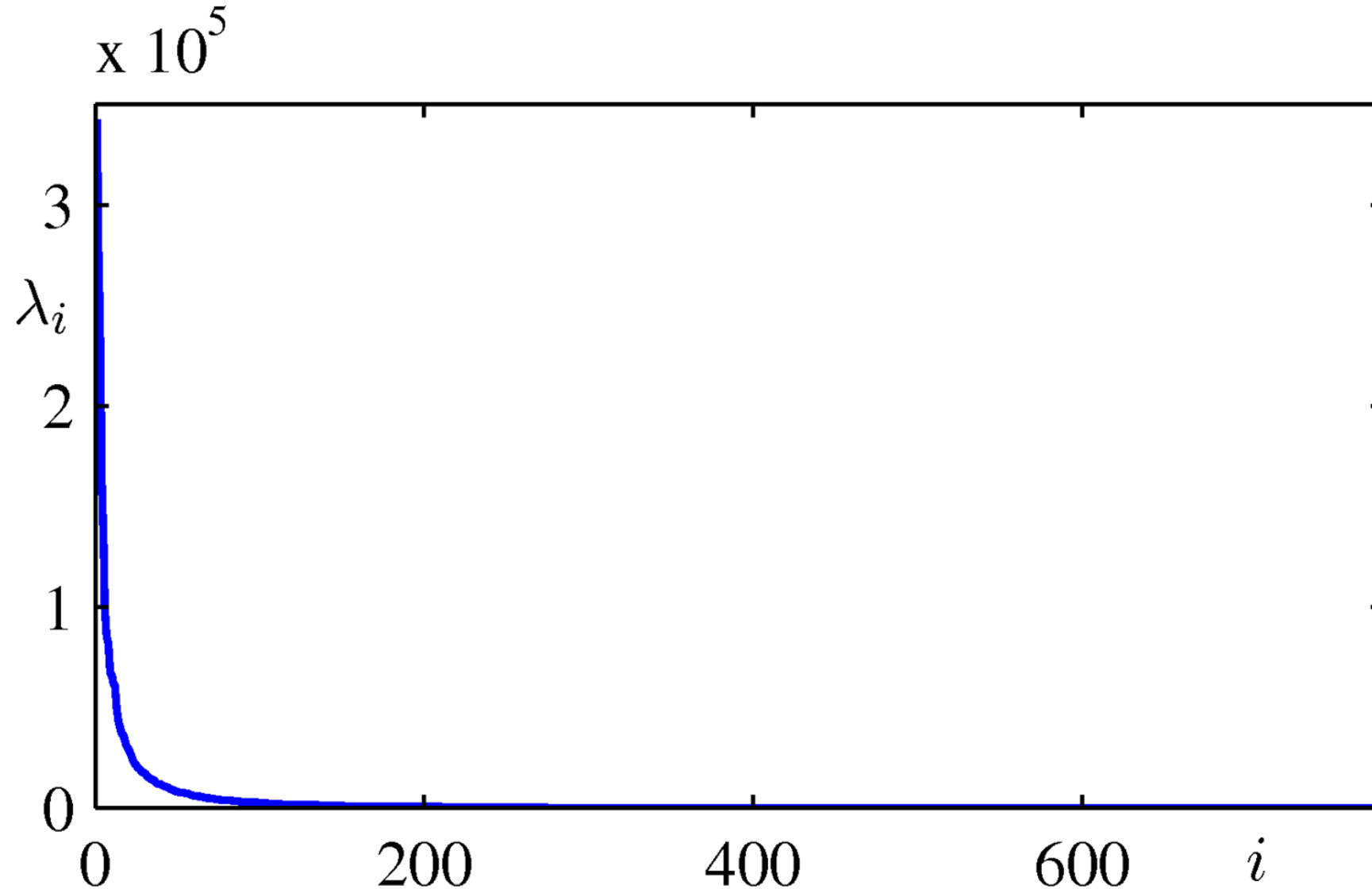
The **fraction of variance explained** =
$$\frac{\sum_{i=1}^M \lambda_i}{\sum_{i=1}^D \lambda_i}$$

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

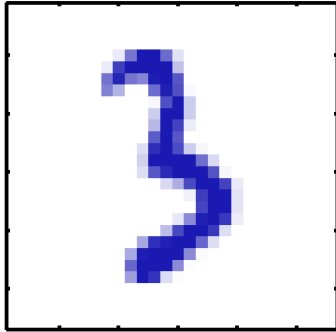
Eigenvalues by principal component i



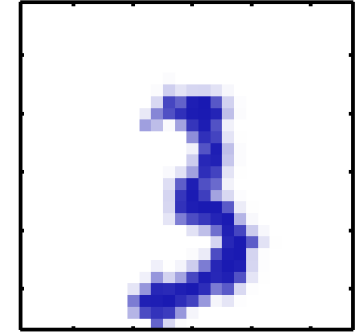
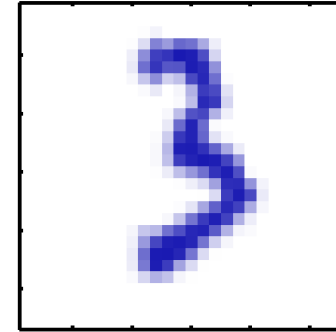
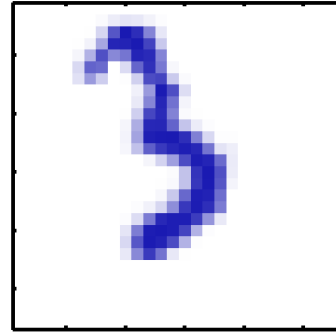
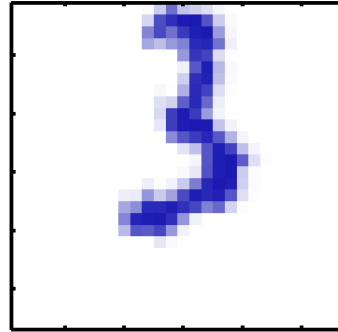
Bishop, Pattern Recognition, 2006

Example: translated digits

Original digit



Translated digits



Types of translation:

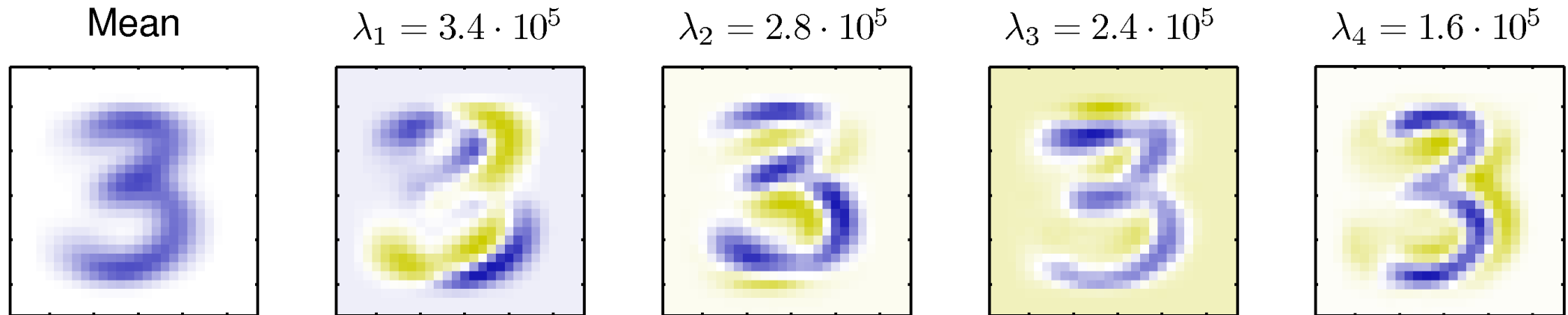
1. Horizontal translation
2. Vertical translation
3. Rotation

Original digits: 64 x 64 pixels

New size: 100 x 100 pixels

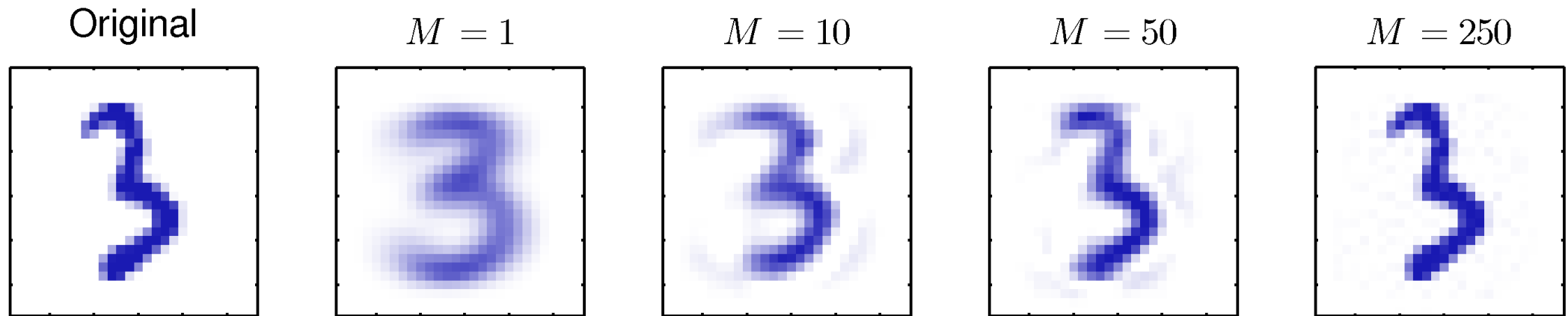
Example: translated digits

Examples of first four principle component eigenvectors and eigenvalues:



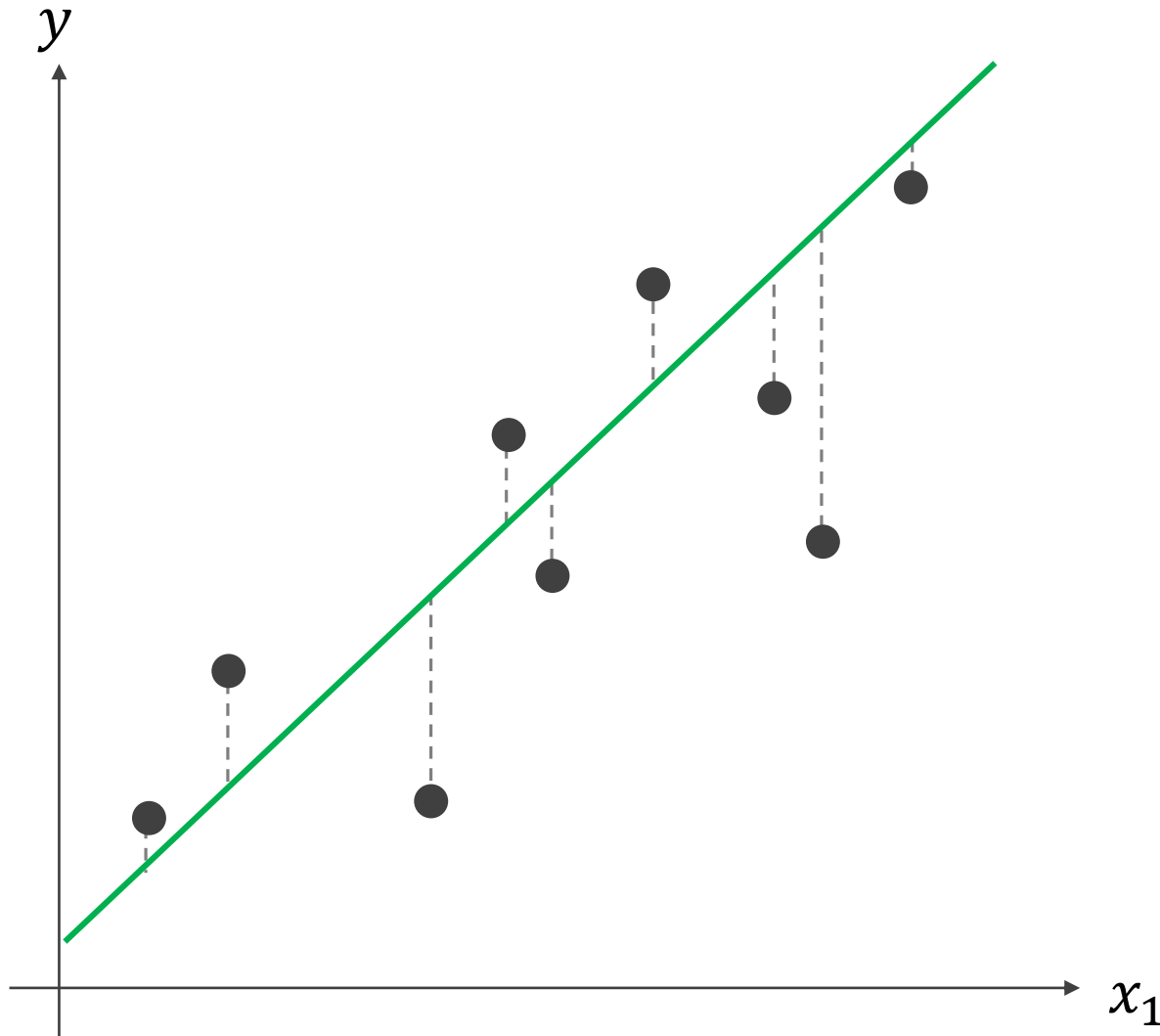
Example: translated digits

Reconstructed examples using different numbers of principal components:

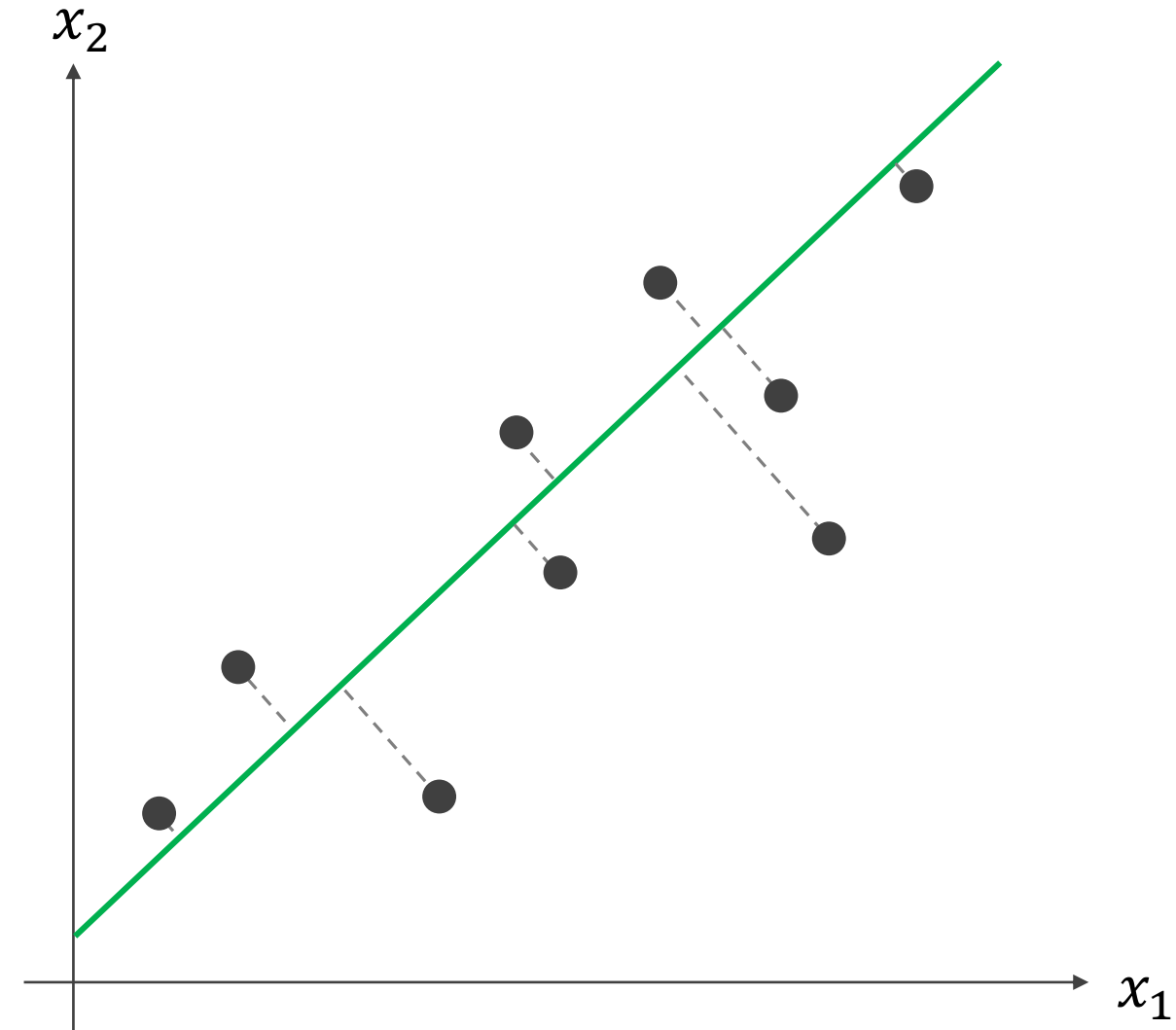


Relationship between the objective for least squares and PCA

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)

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

rows = observations (N)

size $[N \times D]$

Each observation as a vector:

$$\mathbf{x}_i \quad i = 1, \dots, D$$

size $[D \times 1]$

$$\mathbf{u}_i$$

eigenvectors / principal components

size $[D \times 1]$

$$\lambda_i$$

eigenvalues (how much of the variance is explained)

size [scalar]

$i = 1, \dots, D$

$$z_{ij} = \mathbf{u}_j^T \mathbf{x}_i$$

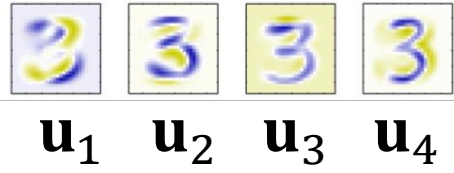
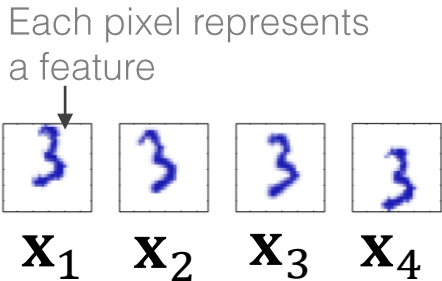
$j = 1, \dots, D$
 $i = 1, \dots, N$

size [scalar]

$$\mathbf{A} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_M]$$

size $[D \times M]$

$$\mathbf{z}_i = \mathbf{A}^T \mathbf{x}_i \quad i = 1, \dots, N$$



$$\mathbf{u}_1 \cdot \mathbf{x}_1 = z_{11}$$

$\mathbf{u}_1 \cdot \mathbf{x}_1 = 2.43$

size $[D \times 1] \quad [D \times 1]$ [scalar]

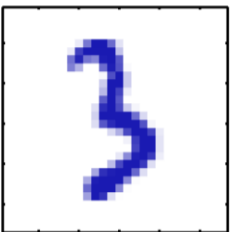
Images from Bishop, Pattern Recognition, 2006

Reconstructing our data from principal components

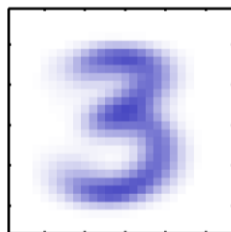
Sum the product of our projected data, \mathbf{z}_i , and our principle components

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

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



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



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

$M = 1$

$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}} + z_{i1} \cdot \mathbf{u}_1 = \text{[blurred 3]} + z_{i1} \cdot \text{[yellow 3]} = \text{[blurred 3]}$$

$M = 250$

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

$M = 10,000$

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

(perfect reconstruction)

Images from Bishop, Pattern Recognition, 2006

Why PCA?

- Dimensionality reduction
- Feature extraction
- Data visualization
- Reducing feature correlation
- Lossy data compression

Other dimensionality reduction techniques

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

e.g. Manifold Learning

