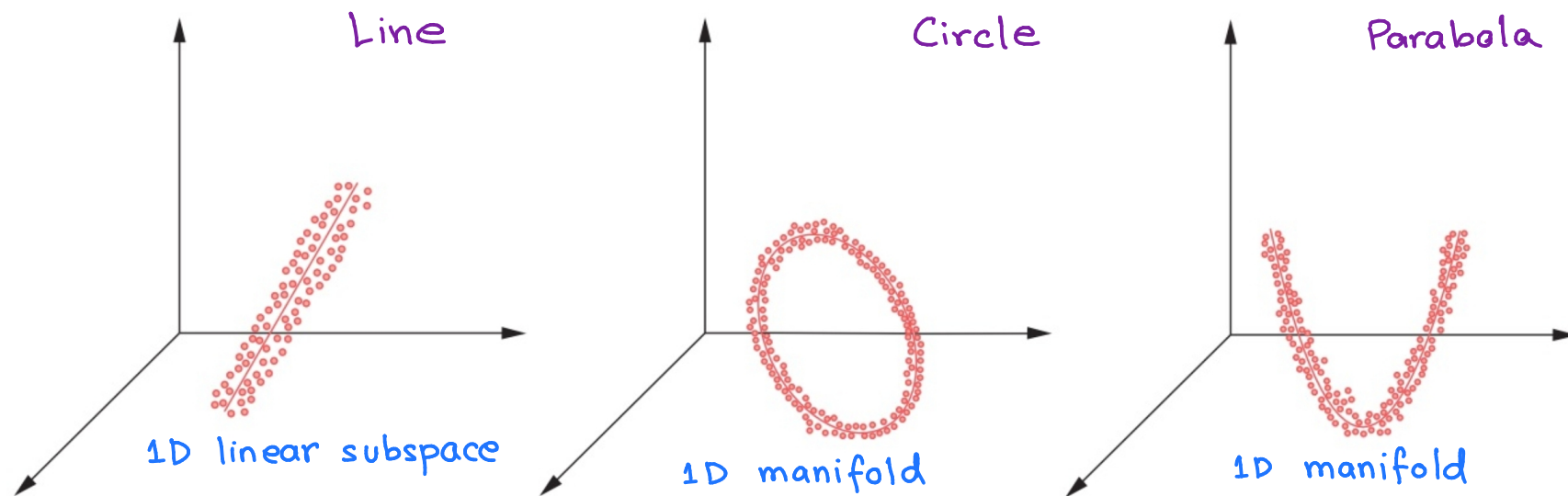


Dimensionality Reduction

- In unsupervised learning, we have seen clustering.
- In this lecture, we will look at dimensionality reduction
- In many practical applications, the input data $\underline{x} \in \mathbb{R}^p$ is a very high-dimensional, however, the **intrinsic dimensionality** may be quite small



In all three cases, the intrinsic dimensionality of data is 1

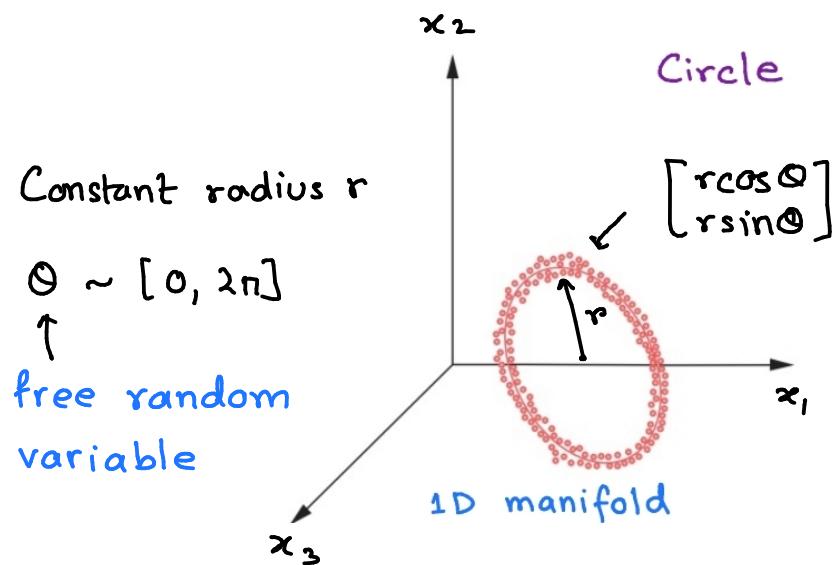
Intrinsic Dimensionality

- A data set $\{\underline{x}^{(i)}\}_{i=1}^N$, with $\underline{x} \in \mathbb{R}^P$, is said to have intrinsic dimensionality $M \leq P$, if the dataset can be described effectively in terms of 'M' free random variables

$$\underline{x} = g(\underline{u})$$

$\mathbb{R}^P \leftarrow \quad \rightarrow \mathbb{R}^M$

Example



The data lies along the circumference of a circle of radius r and a single free parameter Θ suffices to describe the data

Intrinsic dimension = 1

Intrinsic Dimensionality

- An important concern in ML is learning from high-dimensional data \mathbf{x}
- Success of ML, in particular deep learning, is due to its capability of

learning a useful representation of high-dimensional data

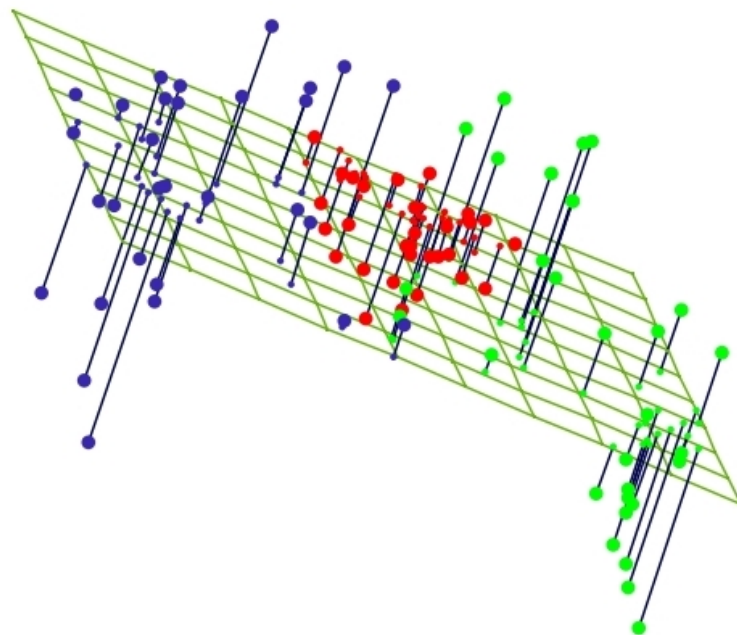
- One of the goals of unsupervised learning:

Learning a lower-dimensional subspace for encoding high-dimensional data set

- Idea of dimensionality reduction: Map data to a lower dimensional space
 - Save computational time in modelling high-dimensional data
 - Visualization in 2-dimensions can offer insights
 - Reduce overfitting and achieve better generalization

Linear Dimensionality Reduction

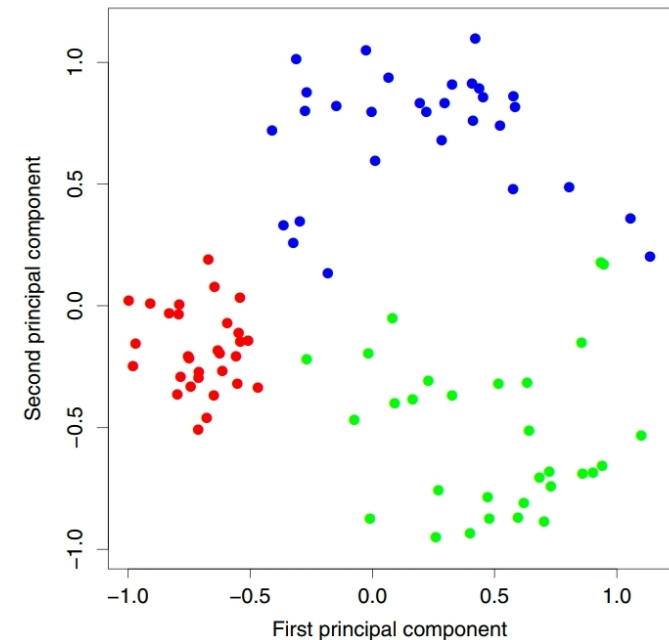
- We will introduce linear dimensionality reduction using Principal Component Analysis (PCA)
- PCA is also known as Karhunen-Loève (KL) transform
 - It falls under linear dimensionality reduction techniques



3D space

Projection
on a linear
subspace

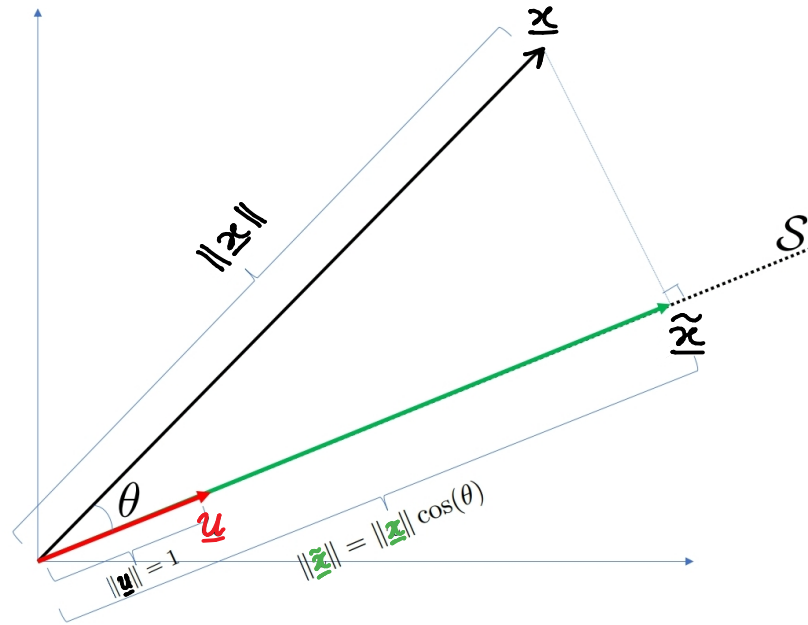
PCA



2D - space

Idea of projection

- Consider projection onto 1-D subspace (a line)



- Subspace S is the line along the unit vector \underline{u}
 - \underline{u} is the basis of S : Any point in S can be written as $z \underline{u}$ for some scalar z

- Projection of vector \underline{x} on S is denoted by $\tilde{\underline{x}} = \text{Proj}_S(\underline{x})$

- Recall that: $\underline{x}^T \underline{u} = \|\underline{x}\| \|\underline{u}\| \cos(\theta) = \|\underline{x}\| \cos \theta = \|\tilde{\underline{x}}\|$

- $\tilde{\underline{x}} = \text{Proj}_S(\underline{x}) = \underbrace{\underline{x}^T \underline{u}}_{\text{length of projection}} \cdot \underbrace{\underline{u}}_{\text{direction of projection}} = \|\tilde{\underline{x}}\| \underline{u}$

Idea of projection

- How to project onto a K -dimensional subspace?
 - **Idea**: Choose an orthonormal bases $\{\underline{u}_1, \underline{u}_2, \dots, \underline{u}_K\}$ for S
 - Project onto each unit vector individually (as in previous slide) and sum together the projections
- Mathematically, the projection is given as:

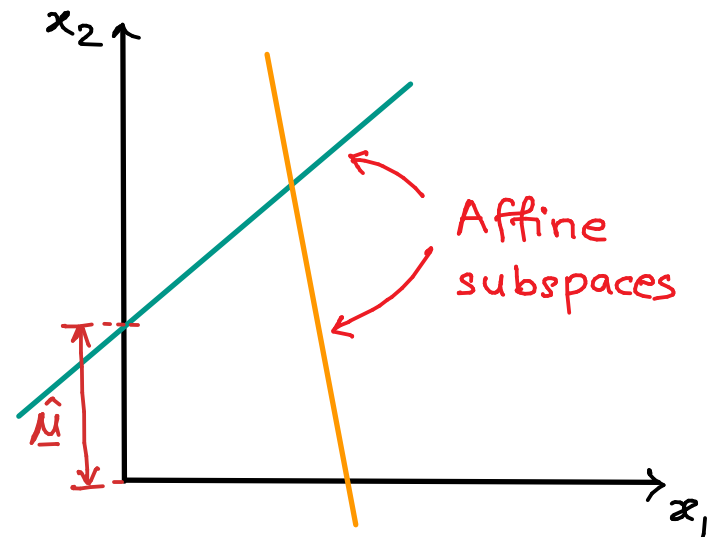
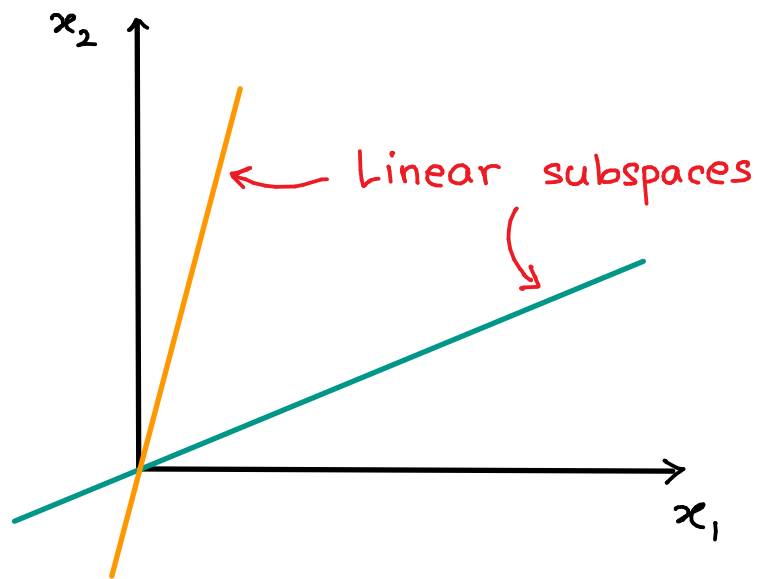
$$\tilde{\underline{x}} = \text{Proj}_S(\underline{x}) = \sum_{i=1}^K z_i \underline{u}_i \quad \text{where} \quad z_i = \underline{x}^T \underline{u}_i$$

- In vector form:

$$\tilde{\underline{x}} = \text{Proj}_S(\underline{x}) = \underline{U} \underline{z} = \begin{bmatrix} | & | & & | \\ \underline{u}_1 & \underline{u}_2 & \dots & \underline{u}_K \\ | & | & & | \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_K \end{bmatrix}, \quad \text{where} \quad \underline{z} = \underline{U}^T \underline{x}$$

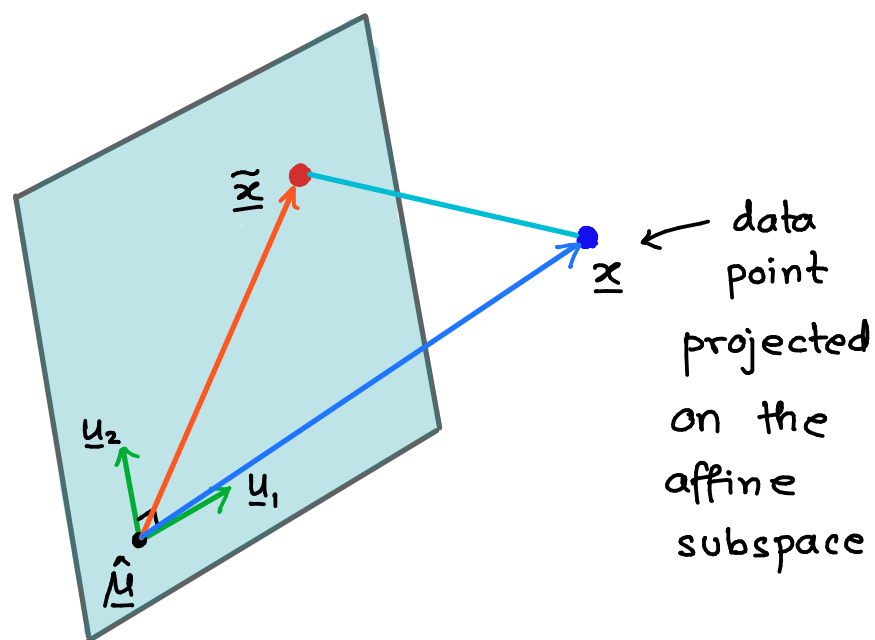
Projection onto an affine subspace

- So far, we have assumed a subspace that passes through zero
- However, the subspaces that we want to project onto can also be **affine subspaces**, which need not pass through zero



The affine subspaces can have an arbitrary origin $\hat{\mu}$

Projection onto an affine subspace



$$\begin{aligned}\tilde{\underline{x}} &= \text{Proj}_S(\underline{x}) \\ &= \underline{U} \underline{z} + \hat{\underline{\mu}} \\ &= z_1 \underline{u}_1 + z_2 \underline{u}_2 + \hat{\underline{\mu}}\end{aligned}$$
$$\underline{z} = \underline{U}^T (\underline{x} - \hat{\underline{\mu}})$$

The affine subspace
has an origin $\hat{\underline{\mu}}$

- $\tilde{\underline{x}}$ is called the **reconstruction** of \underline{x}
- \underline{z} is its **feature / code**
- If all the data points \underline{x} lie close to the subspace, we could approximate \underline{x} with its reconstructions $\tilde{\underline{x}}$

$$\underline{x} \approx \underline{U} \underline{z} + \hat{\underline{\mu}}$$

How to choose a good subspace?

- We want to choose a subspace S which is low-dimensional compared to the dimension of the input space
- How to choose such a subspace S ?
 - We need to find appropriate $\hat{\underline{\mu}}$ and the orthogonal bases $\underline{\underline{U}}$
 - origin $\hat{\underline{\mu}}$ can be set equal to the mean of the dataset
- To find $\underline{\underline{U}}$, one of the two equivalent criteria could be followed:
 - Minimize the reconstruction error:

$$\arg \min_{\underline{\underline{U}}} \frac{1}{N} \sum_{i=1}^N \left\| \underline{x}^{(i)} - \tilde{\underline{x}}^{(i)} \right\|_2^2$$

- Maximize the variance of reconstructions: Find a subspace where the data has the most variability

$$\arg \max_{\underline{\underline{U}}} \frac{1}{N} \sum_{i=1}^N \left\| \underline{x}^{(i)} - \tilde{\underline{x}}^{(i)} \right\|_2^2$$

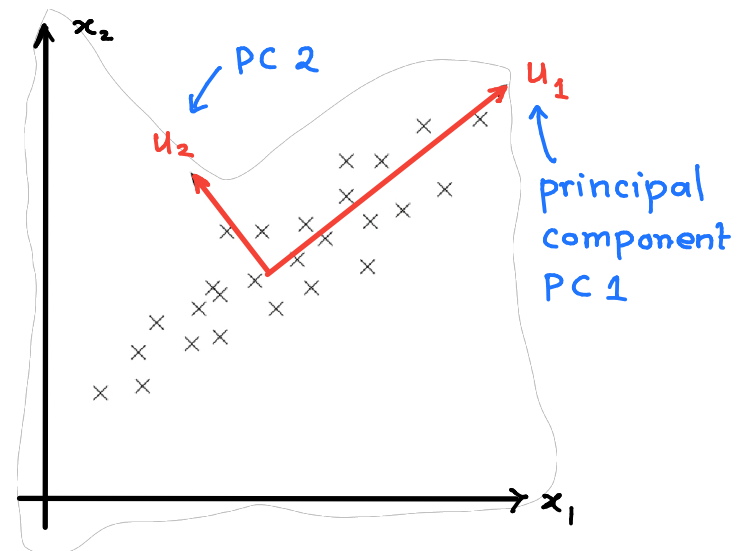
Principal Component Analysis

- Choosing a subspace to maximize the projected variance, or minimize the reconstruction error, is called PCA

- Consider the sample covariance matrix:

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

- $\hat{\Sigma}$ is symmetric and Positive semi-definite (PSD)
- The optimal PCA subspace is spanned by the top 'M' eigenvectors of $\hat{\Sigma}$
- These eigenvectors are called **principal components** or principal directions, much like the principal axes of an ellipse

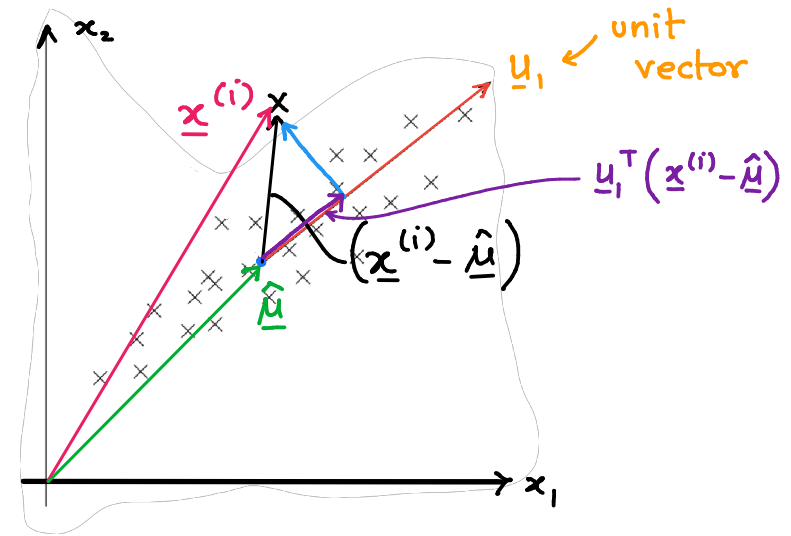


Derivation of PCA

- Let us consider the simplest case of finding a 1-D subspace
 - The goal then is to find a single direction represented by unit vector \underline{u}_1

- Lets maximize the projected variance

$$\begin{aligned} J(\underline{u}_1) &= \frac{1}{N} \sum_{i=1}^N \left(\underline{u}_1^T (\underline{x}^{(i)} - \hat{\underline{\mu}}) \right)^2 \\ &= \frac{1}{N} \sum_{i=1}^N \underline{u}_1^T (\underline{x}^{(i)} - \hat{\underline{\mu}}) (\underline{x}^{(i)} - \hat{\underline{\mu}})^T \underline{u}_1 \\ &= \underline{u}_1^T \hat{\sum} \underline{u}_1 \end{aligned}$$



- So the optimization task is:

$$\begin{aligned} \underline{u}_1 &= \underset{\underline{u}}{\operatorname{argmax}} \quad \underline{u}^T \hat{\sum} \underline{u} \\ \text{s.t.} \quad &\underline{u}^T \underline{u} = 1 \end{aligned}$$

Lagrangian: $L(\underline{u}, \lambda) = \underline{u}^T \hat{\sum} \underline{u} - \lambda (\underline{u}^T \underline{u} - 1)$

Take gradient and set to zero:

$$\hat{\sum} \underline{u} = \lambda \underline{u}$$

eigenvalue
eigenvector

\therefore Principal direction \underline{u}_1 is an eigenvector

- Since $\hat{\Sigma}$ is symmetric and PSD, all eigenvalues are **real** and **non-negative**: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$
- The 2nd principal component \underline{u}_2 is selected such that:
 - (a) \underline{u}_2 is orthogonal to \underline{u}_1
 - (b) \underline{u}_2 maximizes the variance after projecting the data onto the direction of \underline{u}_2
 - (c) The **2nd principal component** (or direction) is the eigenvector corresponding to the **2nd largest eigenvalue** of $\hat{\Sigma}$, λ_2
- Similar arguments can be used to show that the '**m**'th **principal component** is the '**m**'th eigenvector of $\hat{\Sigma}$
- The process continues until **M** principal components (corresponding to the **M** largest eigenvalues)

PCA decorrelates features

- The features (or code) are decorrelated by PCA

$$\text{Cov}(\underline{z}) = \text{Cov}(\underline{U}^T (\underline{x} - \hat{\underline{\mu}}))$$

$$= \underline{U}^T \text{Cov}(\underline{x}) \underline{U}$$

$$= \underline{U}^T \hat{\underline{\Sigma}} \underline{U}$$

$$= \underline{U}^T \underline{Q} \underline{\Lambda} \underline{Q}^T \underline{U}$$

$$= \begin{bmatrix} \underline{I} & \underline{0} \end{bmatrix} \underline{\Lambda} \begin{bmatrix} \underline{I} \\ \underline{0} \end{bmatrix}$$

$$= \text{top left } M \times M \text{ block} \\ \text{of } \underline{\Lambda}$$

Spectral decomposition

$$\hat{\underline{\Sigma}}_{P \times P} = \underline{Q} \underline{\Lambda} \underline{Q}^T$$

\underline{Q} is the eigenvector matrix (orthonormal)
 $\underline{\Lambda}$ is the eigenvalues matrix

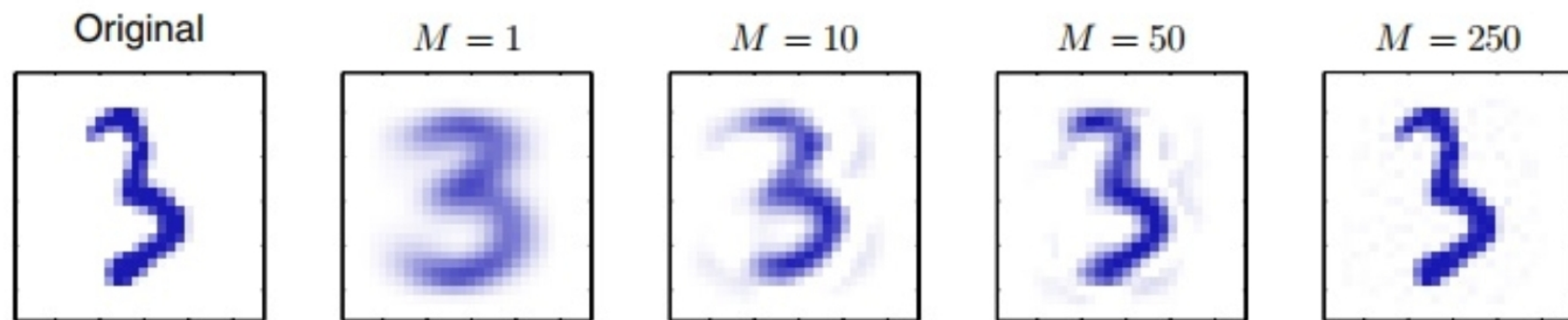
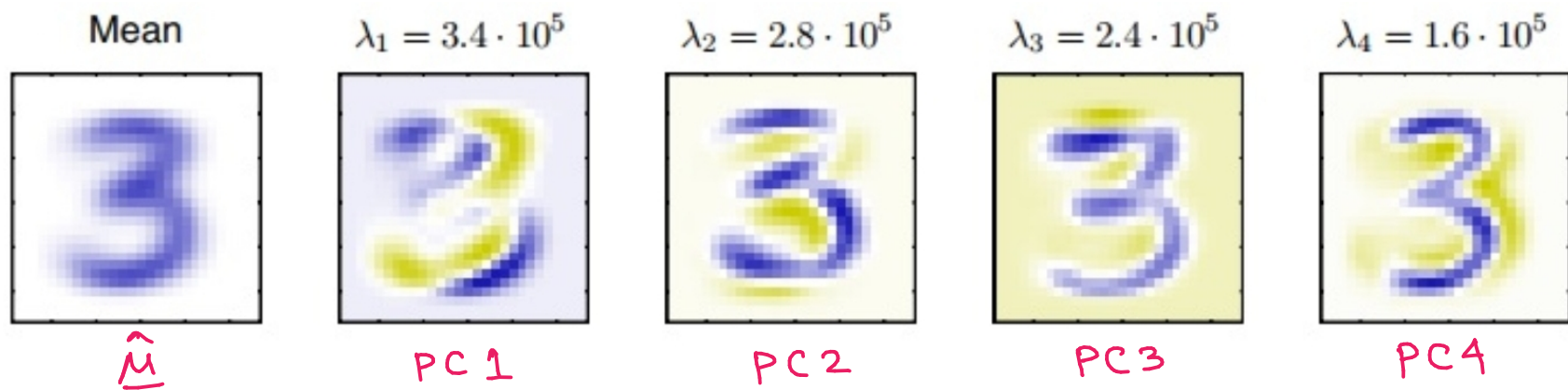
$$\underline{Q}_{P \times P} = \begin{bmatrix} \underline{U}_{P \times M} & \vdots & \underline{U}_{\perp, P \times (P-M)} \end{bmatrix}$$

- Covariance of feature \underline{z} is diagonal \rightarrow uncorrelated

Summary of PCA

- Dimensionality reduction aims to find a low-dimensional representation of the data
- PCA projects the data onto an affine subspace that maximizes projected variance or minimizes the reconstruction error
- The optimal subspace is given by the top M eigenvectors of the sample covariance matrix, corresponding to the M largest eigenvalues
- PCA gives a set of decorrelated features

Example of data compression



Original
digit

PCA reconstructions