

# Introduction to Machine Learning:

## Lecture 4 – Unsupervised Learning

Michael Kagan



TRISEP Summer School  
July 8-12, 2024

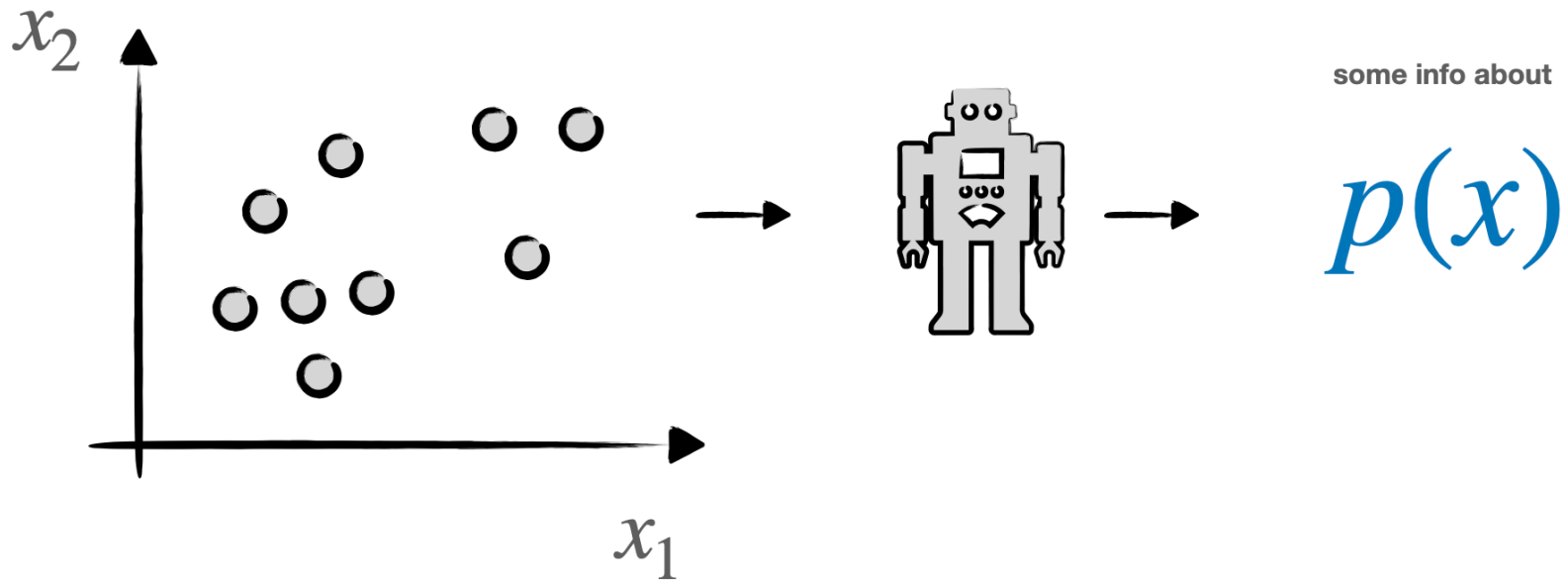
- Lecture 1 – Machine Learning Fundamentals
- Lecture 2 – Intro to Neural Networks
- Lecture 3 – Intro to Deep Learning
- Lecture 4 – Intro to Unsupervised Learning
- Lecture 5 – Intro to Deep Generative Models

---

# Beyond Regression and Classification

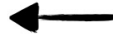
- Not all tasks are predicting a label from features, as in classification and regression
- May want to model a high-dim. signal
  - Data synthesis / simulation
  - Density estimation
  - Anomaly detection
  - Denoising, super resolution
  - Data compression
  - ...
- Often don't have labels → Unsupervised Learning

- Our goal is to study the data density  $p(x)$
- Even w/o labels, aim to characterize the distribution



**A process**

$$\text{Die} \rightarrow \mathbb{R}^2$$



**A formula**

$$\mathbb{R}^2 \rightarrow \mathbb{R}$$

$$p_{\mu, \Sigma}(x) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma (x - \mu)\right)$$

**Evaluating the Probability  
for a given sample**

"Understanding  $p(x)$ " – ability to do either or both of these

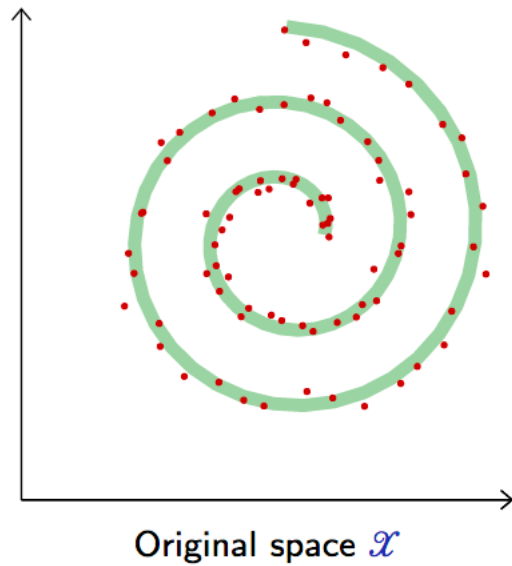
- In many cases, we don't have a theory of the underlying process → *Can still learn to sample*
- Deep learning can be very good at this!

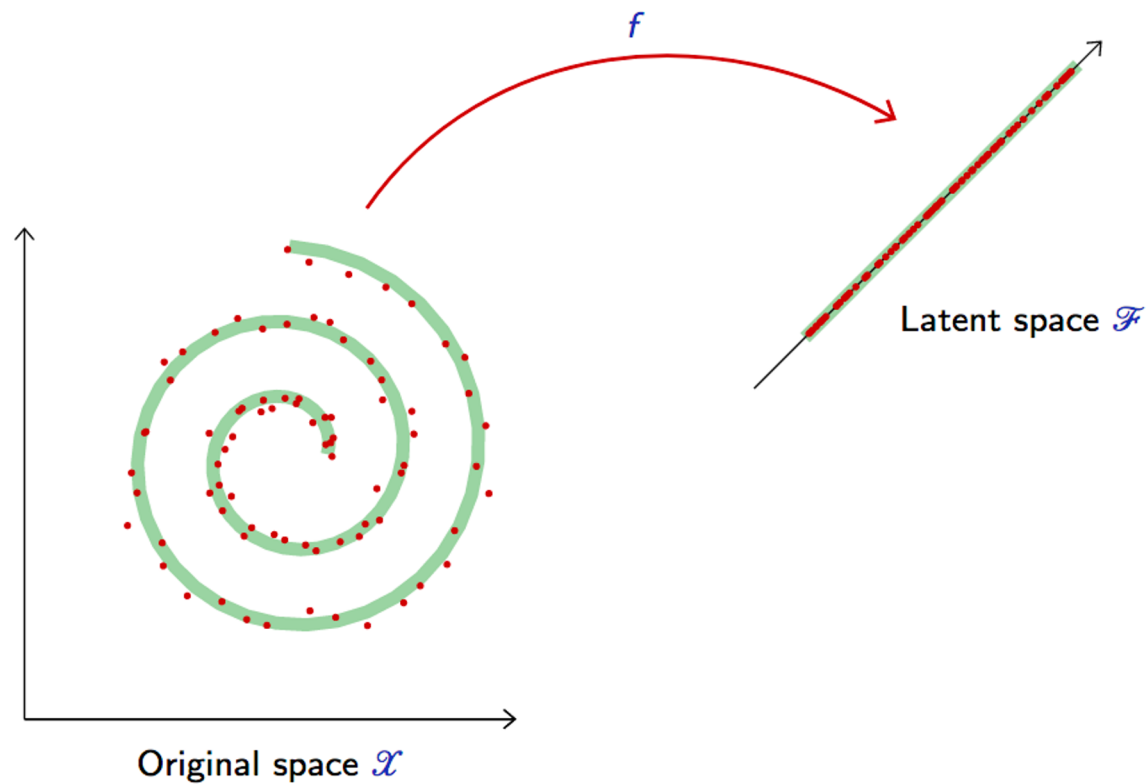
face  $\sim p(\text{face})$

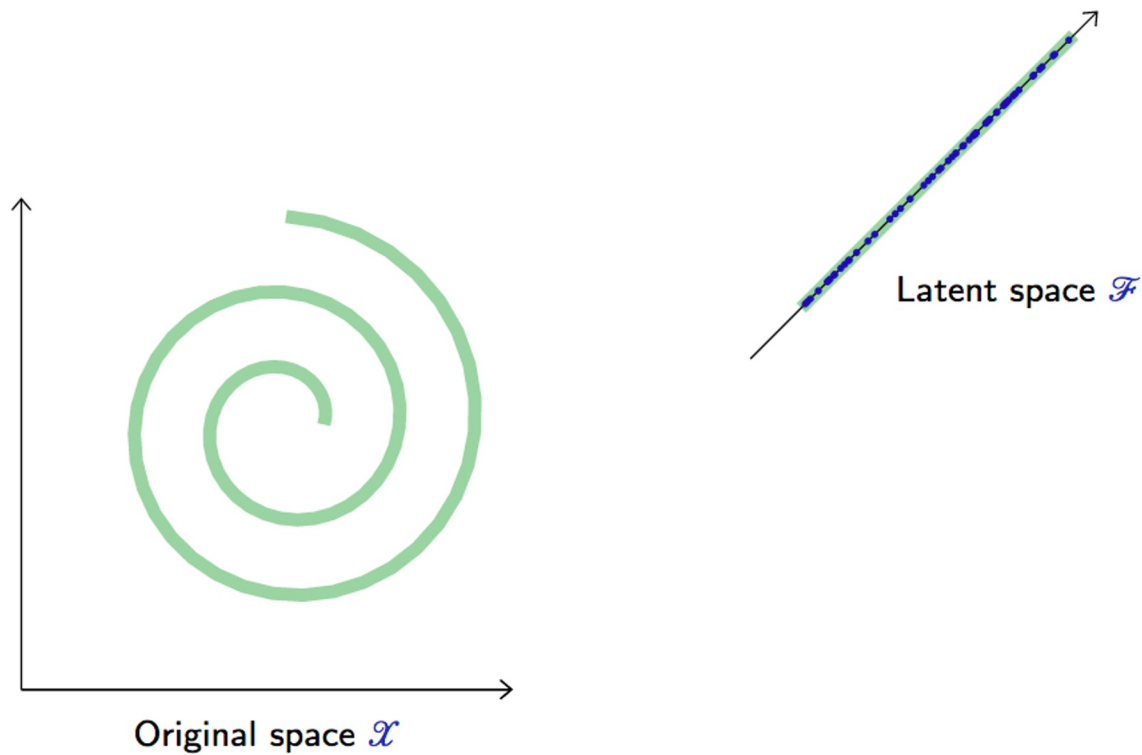


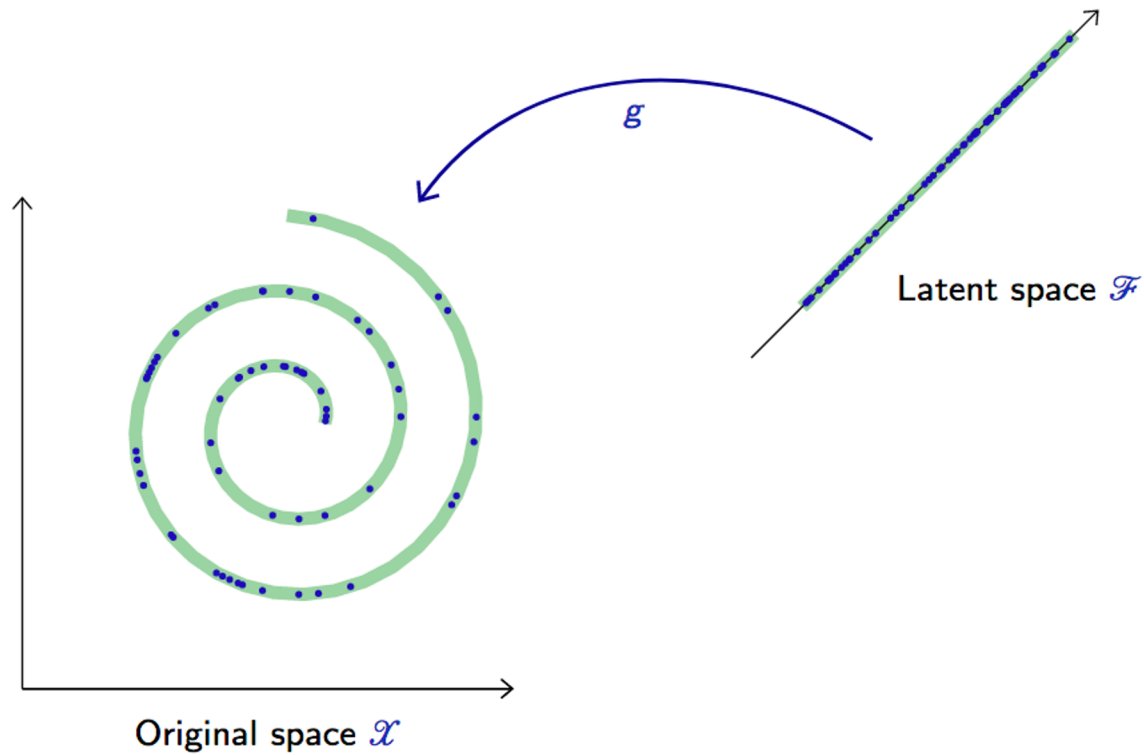
- Unsupervised learning is more heterogeneous than supervised learning
- Many architectures, losses, learning strategies
- Often constructed so model converges to  $p(x)$ 
  - Variational inference, Adversarial learning, Self-supervision, ...
- Often framed as **modeling the lower dimensional “meaningful degrees of freedom”** that describe the data



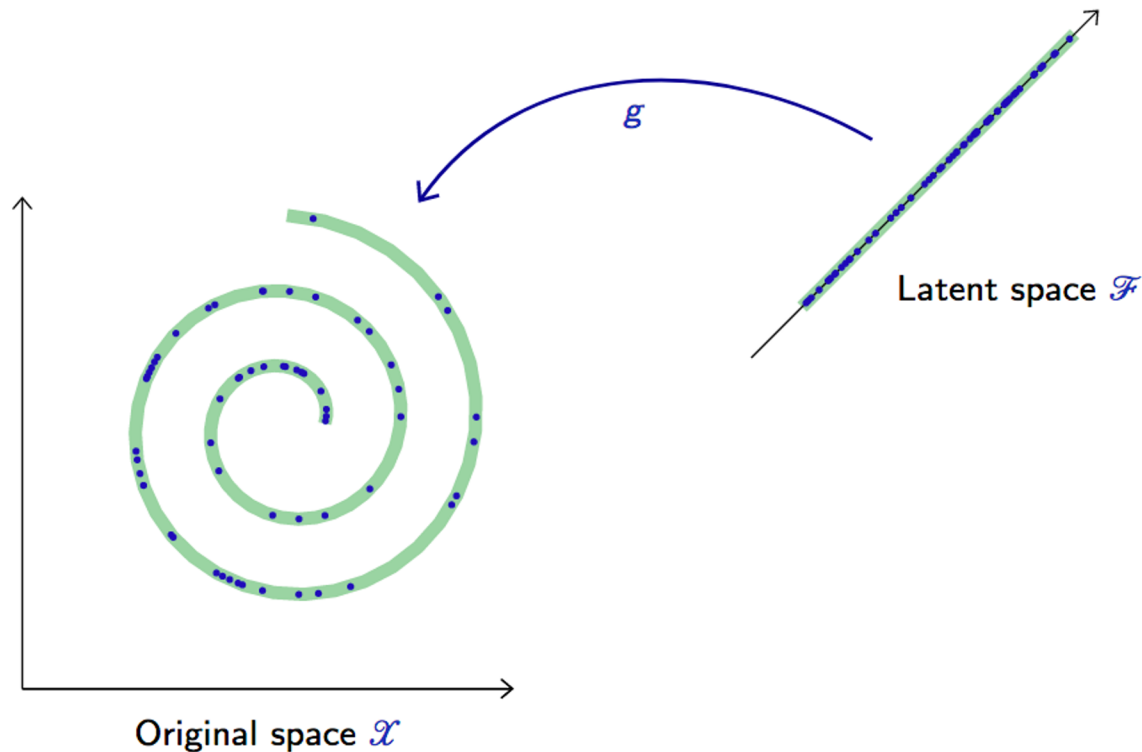








How can we find the “meaningful degrees of freedom” in the data?



- Dimensionality Reduction / Compression
- Can we learn to:
  1. Compress the data to a *latent space* with smaller number of dimensions
  2. Recover the original data from this latent space?
- Latent space must encode and retain the important information about the data

# Principle Components Analysis

Find a low dimensional (less complex) representation of the data with a mapping  $Z=h(X)$



- Given data  $\{\mathbf{x}_i\}_{i=1\dots N}$  can we find a directions in features space that explain most variation of data?

- Given data  $\{\mathbf{x}_i\}_{i=1\dots N}$  can we find a directions in features space that explain most variation of data?
- Data covariance: 
$$\mathbf{S} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^2$$

- Given data  $\{\mathbf{x}_i\}_{i=1\dots N}$  can we find a directions in features space that explain most variation of data?

- Data covariance: 
$$\mathbf{S} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^2$$

- Let  $\mathbf{u}_1$  be the projected direction, we can solve:

$$\mathbf{u}_1^* = \arg \max_{\mathbf{u}_1} \underbrace{\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1}_{\text{Variance of projected data}} + \lambda \underbrace{(1 - \mathbf{u}_1^T \mathbf{u}_1)}_{\text{Unit length vector constraint}}$$

- Given data  $\{\mathbf{x}_i\}_{i=1\dots N}$  can we find a directions in features space that explain most variation of data?

- Data covariance: 
$$\mathbf{S} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^2$$

- Let  $\mathbf{u}_1$  be the projected direction, we can solve:

$$\mathbf{u}_1^* = \arg \max_{\mathbf{u}_1} \underbrace{\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1}_{\text{Variance of projected data}} + \lambda \underbrace{(1 - \mathbf{u}_1^T \mathbf{u}_1)}_{\text{Unit length vector constraint}}$$
$$\nabla_{\mathbf{u}_1} [\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda(1 - \mathbf{u}_1^T \mathbf{u}_1)] = \mathbf{S} \mathbf{u}_1 - \lambda \mathbf{u}_1 = 0$$

- Given data  $\{\mathbf{x}_i\}_{i=1\dots N}$  can we find a directions in features space that explain most variation of data?

- Data covariance: 
$$\mathbf{S} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^2$$

- Let  $\mathbf{u}_1$  be the projected direction, we can solve:

$$\mathbf{u}_1^* = \arg \max_{\mathbf{u}_1} \underbrace{\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1}_{\text{Variance of projected data}} + \lambda \underbrace{(1 - \mathbf{u}_1^T \mathbf{u}_1)}_{\text{Unit length vector constraint}}$$
$$\rightarrow \mathbf{S} \mathbf{u}_1 = \lambda \mathbf{u}_1$$

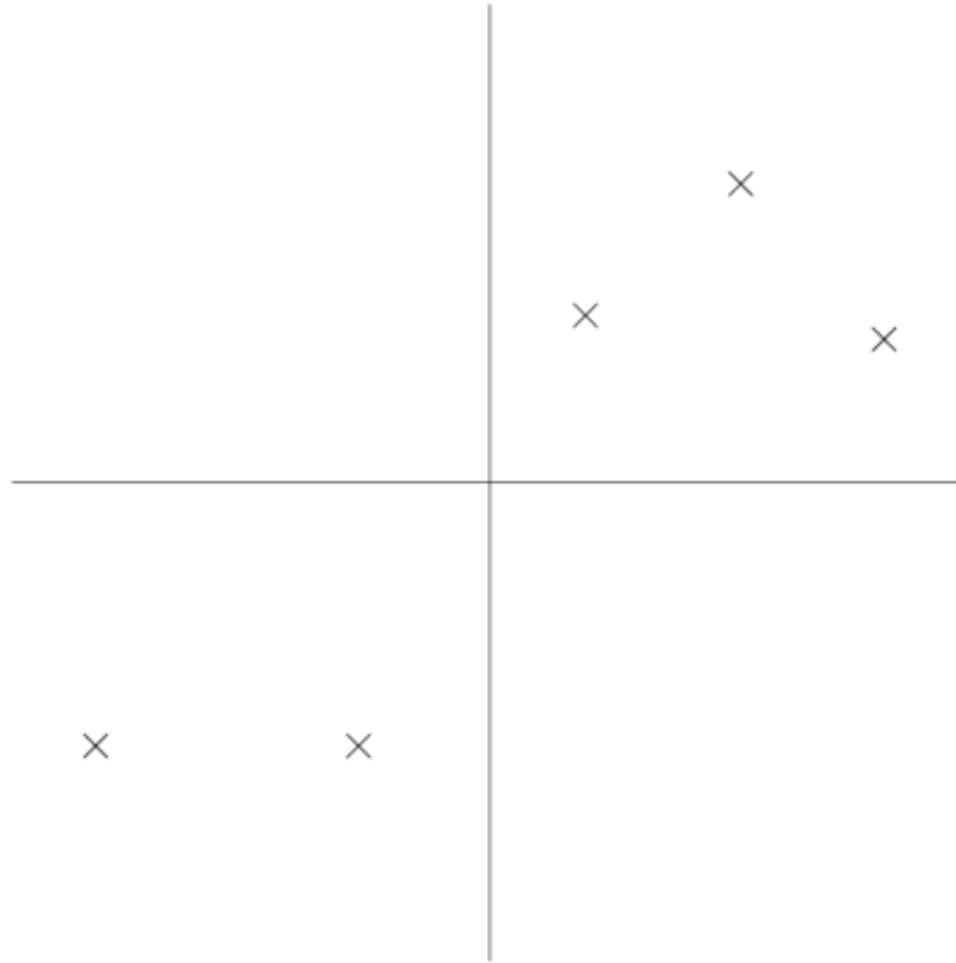
- Given data  $\{\mathbf{x}_i\}_{i=1\dots N}$  can we find a directions in features space that explain most variation of data?

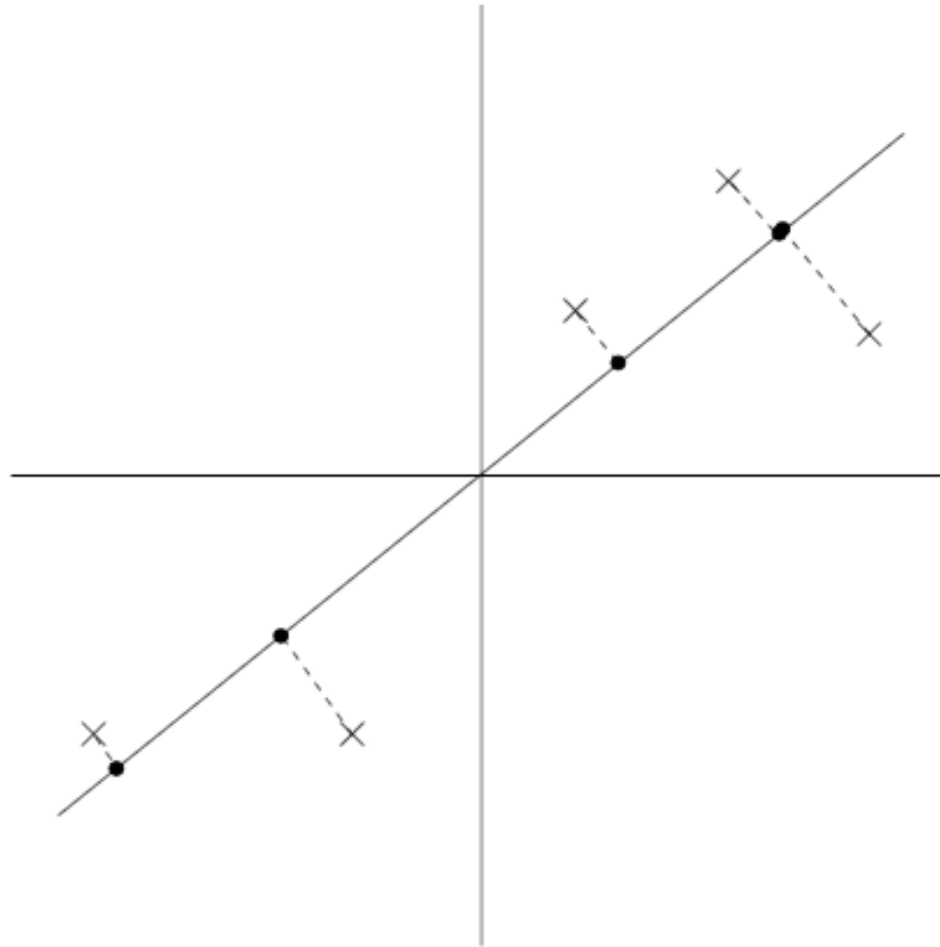
- Data covariance: 
$$\mathbf{S} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^2$$

- Let  $\mathbf{u}_1$  be the projected direction, we can solve:

$$\mathbf{u}_1^* = \arg \max_{\mathbf{u}_1} \underbrace{\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1}_{\text{Variance of projected data}} + \lambda \underbrace{(1 - \mathbf{u}_1^T \mathbf{u}_1)}_{\text{Unit length vector constraint}}$$
$$\rightarrow \mathbf{S} \mathbf{u}_1 = \lambda \mathbf{u}_1$$

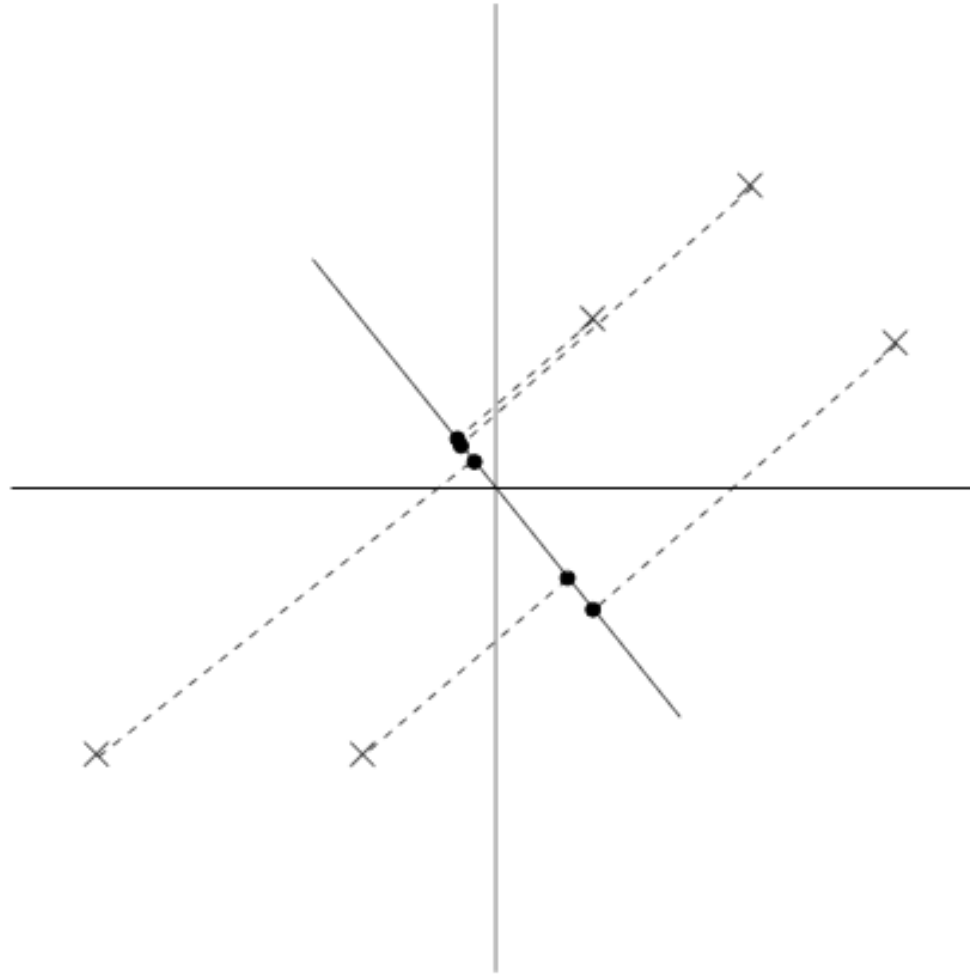
- *Principle components* are the **eigenvectors** of the data covariance matrix!
  - Eigenvalues are the variance explained by that component





First principle component, projects on to this axis have large variance



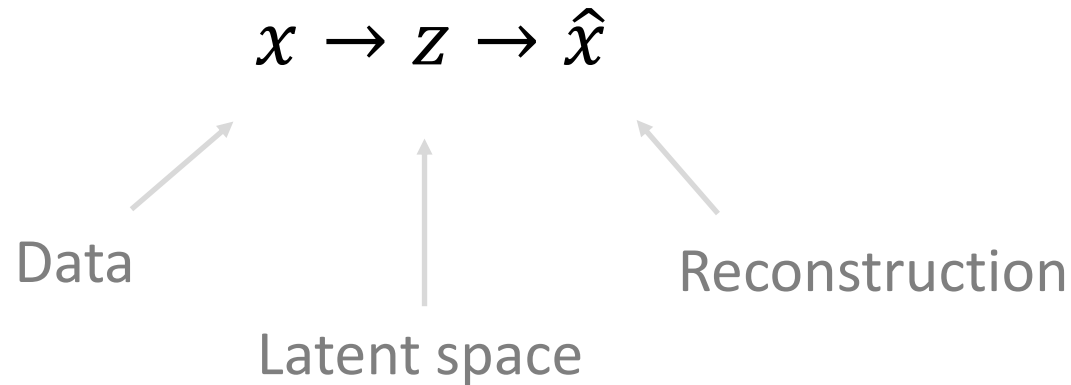


Second principle component, projects have small variance

---

# Autoencoders

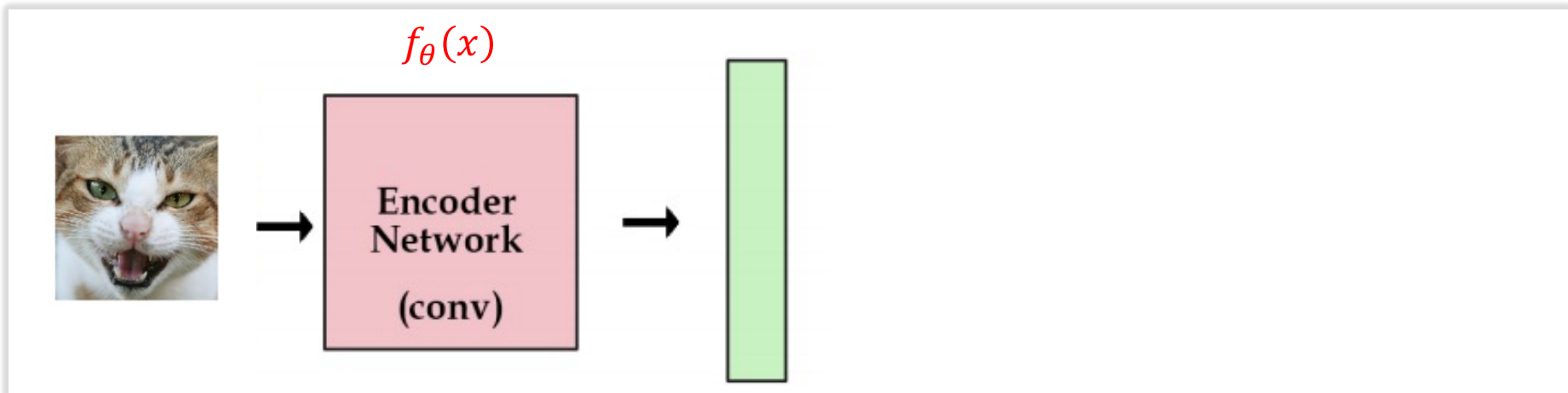
- Map a space to itself through a compression



- Map a space to itself through a compression

$$x \rightarrow z \rightarrow \hat{x}$$

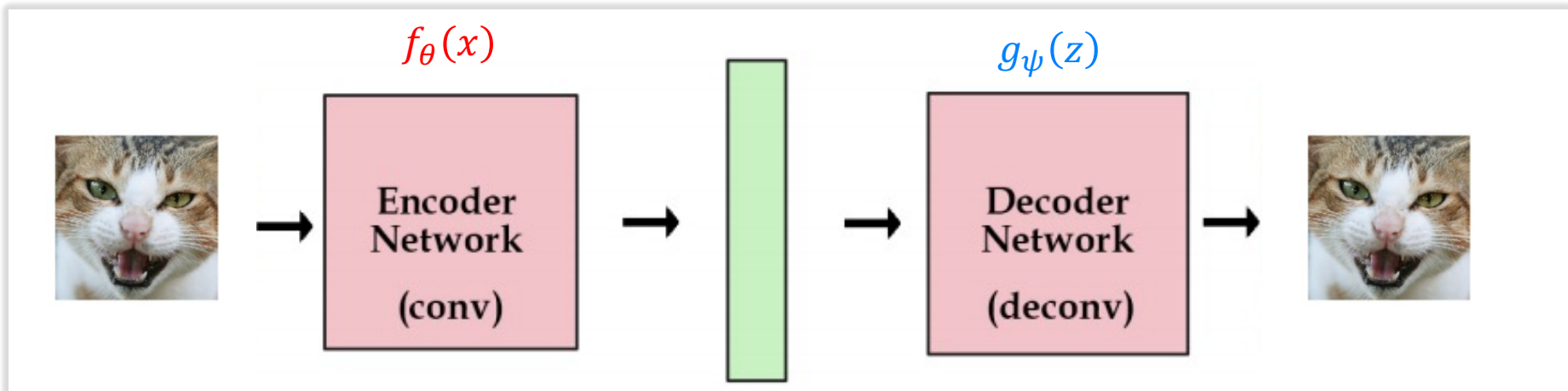
- **Encoder**: Map from data to a lower dim. latent space
  - Neural network  $f_{\theta}(x)$  with parameters  $\theta$

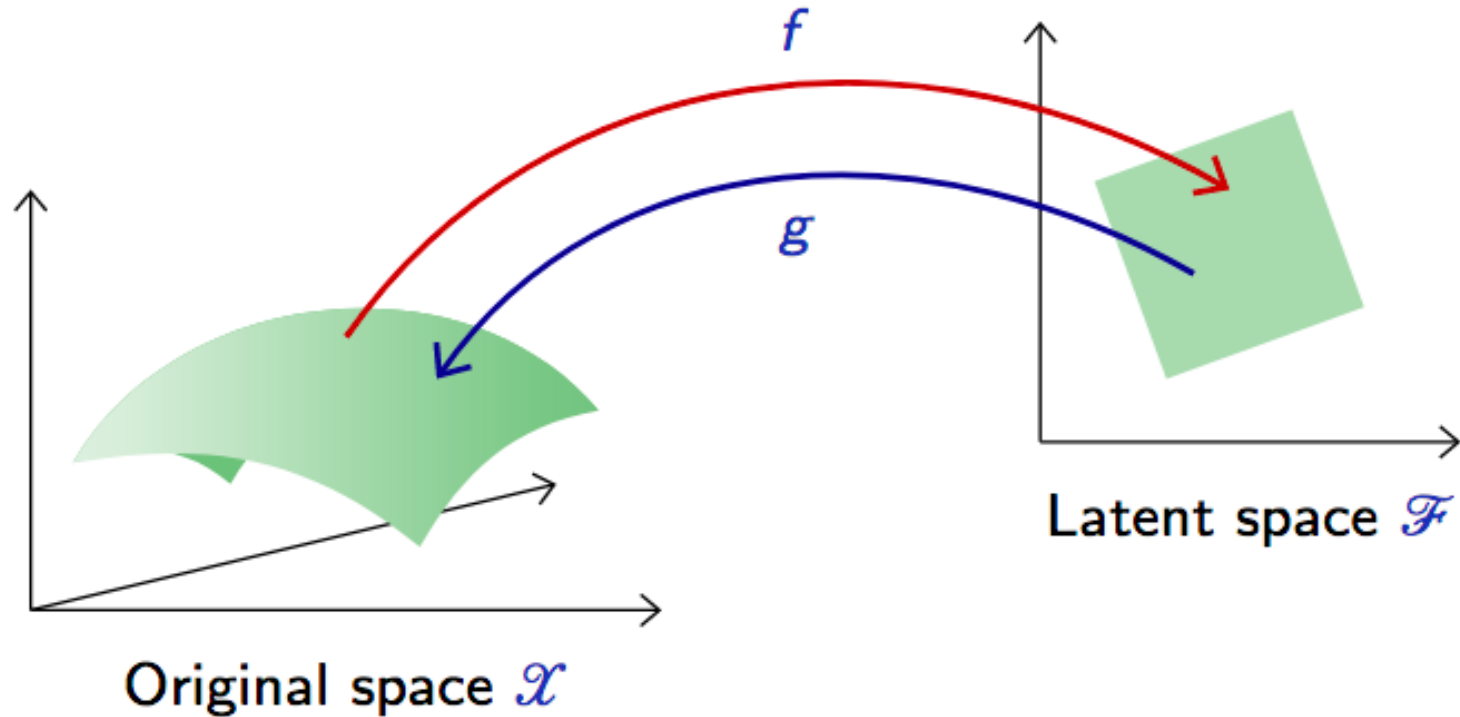


- Map a space to itself through a compression

$$x \rightarrow z \rightarrow \hat{x}$$

- **Encoder**: Map from data to a lower dim. latent space
  - Neural network  $f_{\theta}(x)$  with parameters  $\theta$
- **Decoder**: Map from latent space back to data space
  - Neural network  $g_{\psi}(z)$  with parameters  $\psi$

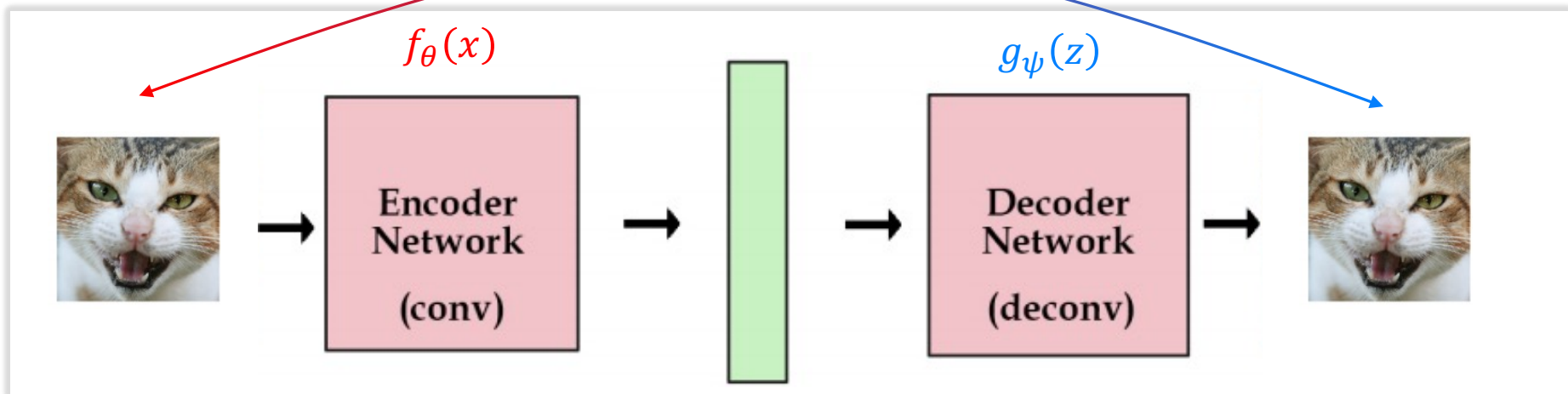




- Latent space is of lower dimension than data
- Model must learn a “good” parametrization and capture dependencies between components

$$L(\theta, \psi) = \frac{1}{N} \sum_n \|x_n - g_\psi(f_\theta(x_n))\|^2$$

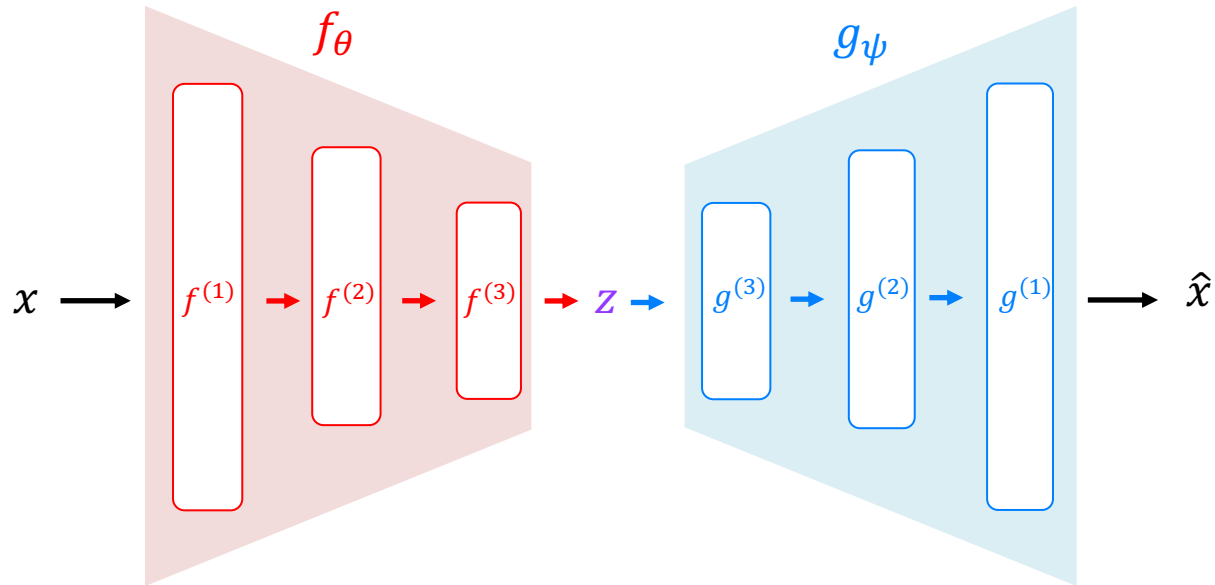
- **Loss:** mean *reconstruction loss* (MSE) between data and encoded-decoded data
- Min. over params. of encoder ( $\theta$ ) and decoder ( $\psi$ ).



$$L(\theta, \psi) = \frac{1}{N} \sum_n \|x_n - g_\psi(f_\theta(x_n))\|^2$$

- Loss: mean *reconstruction loss* (MSE) between data and encoded-decoded data
- Min. over params. of encoder ( $\theta$ ) and decoder ( $\psi$ ).
- NOTE: if  $f_\theta(x)$  and  $g_\psi(z)$  are linear, optimal solution given by Principle Components Analysis





- When  $f_{\theta}$  and  $g_{\psi}$  are multiple neural network layers, can learn complex mappings
  - $f_{\theta}$  and  $g_{\psi}$  can be Fully Connected, CNNs, RNNs, etc.
  - Choice of network structure will depend on data

# Deep Convolutional Autoencoder

$X$  (original samples)

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 8 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

$g \circ f(X)$  (CNN,  $d = 16$ )

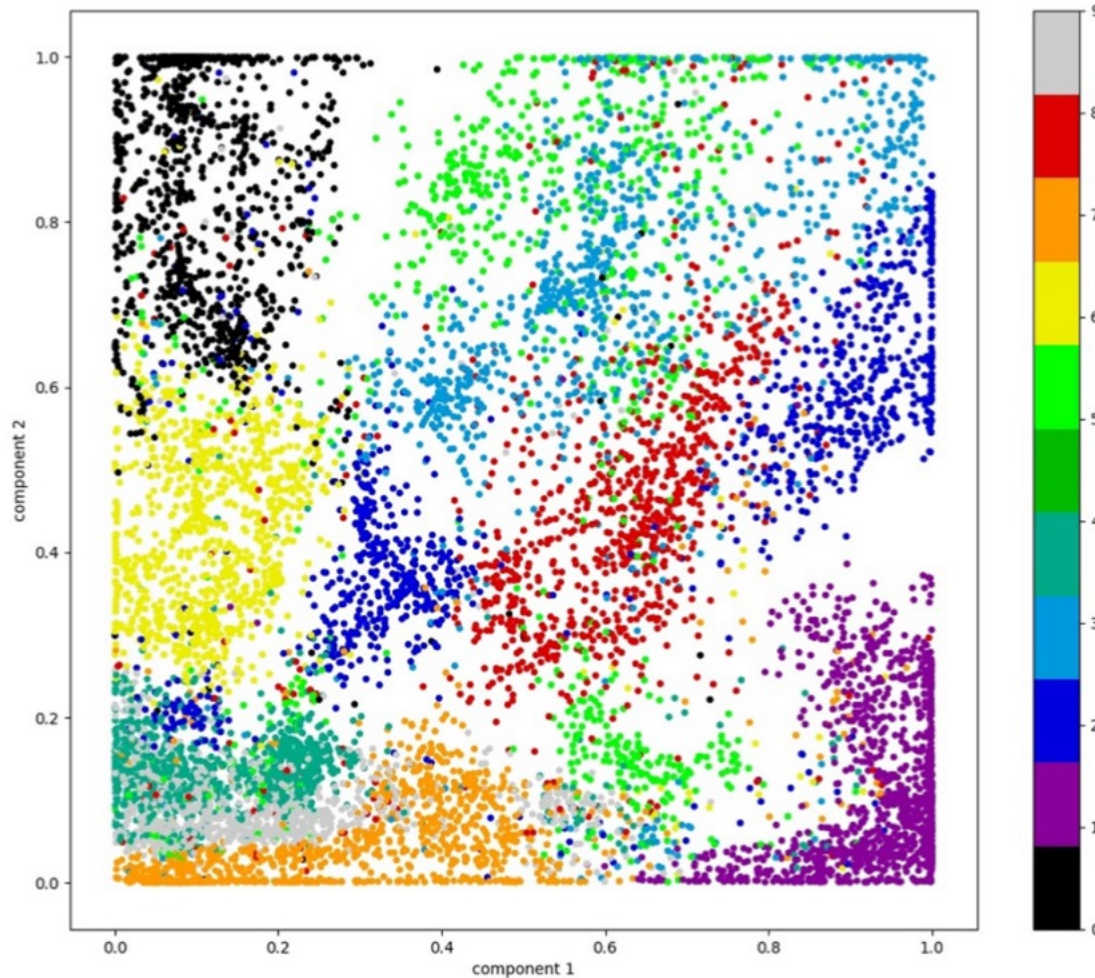
7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 8 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

$g \circ f(X)$  (PCA,  $d = 16$ )

7 2 1 0 9 1 4 9 9 9 0 6  
9 0 1 5 9 7 8 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

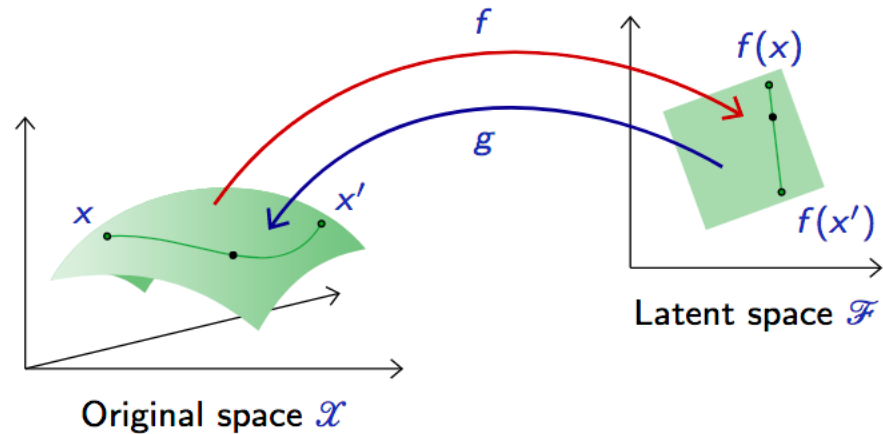
$f_\theta$  and  $g_\psi$  are five convolutional layers

- Can look at latent space to see how the model arranges the data

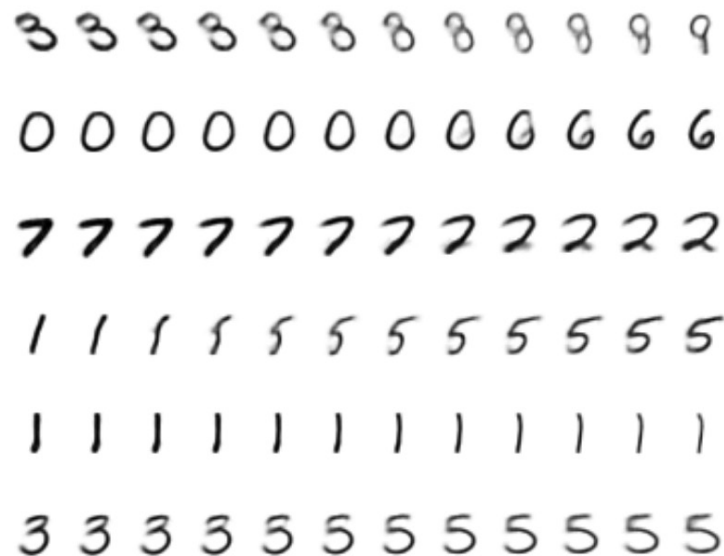


# Interpolating in Latent Space

$$\alpha \in [0, 1], \quad \xi(x, x', \alpha) = g((1 - \alpha)f(x) + \alpha f(x')).$$

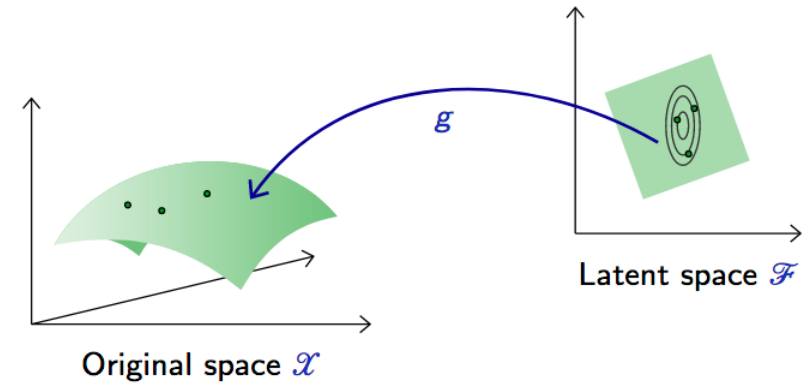


Autoencoder interpolation ( $d = 8$ )



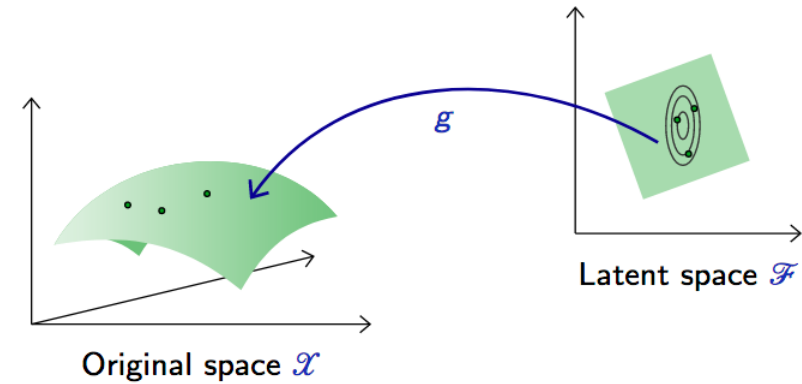
# Can We Generate Data with Decoder?

- Can we sample in latent space and decode to generate data?



# Can We Generate Data with Decoder?

- Can we sample in latent space and decode to generate data?



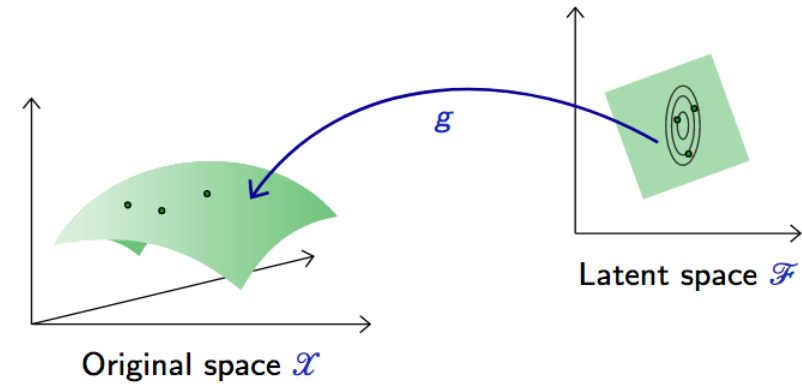
- What distribution to sample from in latent space?
  - Try Gaussian with mean and variance from data

Autoencoder sampling ( $d = 16$ )



# Can We Generate Data with Decoder?

- Can we sample in latent space and decode to generate data?



- What distribution to sample from in latent space?
  - Try Gaussian with mean and variance from data

Autoencoder sampling ( $d = 16$ )



- Don't know the right latent space density

- Unsupervised learning aims to characterize data, even without distributions
- Often framed as learning the meaningful degrees of freedom of a system
- We saw examples of powerful ways to learn these meaningful degrees of freedom, e.g. linearly with PCA and non-linearly with autoencoders



---

# Backup