

Autoencoders and Generative Models

CISC 7026: Introduction to Deep Learning

University of Macau

Agenda

1. Review
2. Compression
3. Autoencoders
4. Applications
5. Variational Modeling
6. VAE Implementation
7. Coding

Agenda

1. Review
2. Compression
3. Autoencoders
4. Applications
5. Variational Modeling
6. VAE Implementation
7. Coding

Convolution works over inputs of any variables (time, space, etc)

Review

Convolution works over inputs of any variables (time, space, etc)

Recurrent neural networks only work with time

Review

Convolution works over inputs of any variables (time, space, etc)

Recurrent neural networks only work with time

Convolution makes use of locality and translation equivariance properties

Review

Convolution works over inputs of any variables (time, space, etc)

Recurrent neural networks only work with time

Convolution makes use of locality and translation equivariance properties

Recurrent models do not assume locality or equivariance

Review

Convolution works over inputs of any variables (time, space, etc)

Recurrent neural networks only work with time

Convolution makes use of locality and translation equivariance properties

Recurrent models do not assume locality or equivariance

Equivariance and locality make learning more efficient, but not all problems have this structure

How do humans process temporal data?

How do humans process temporal data?

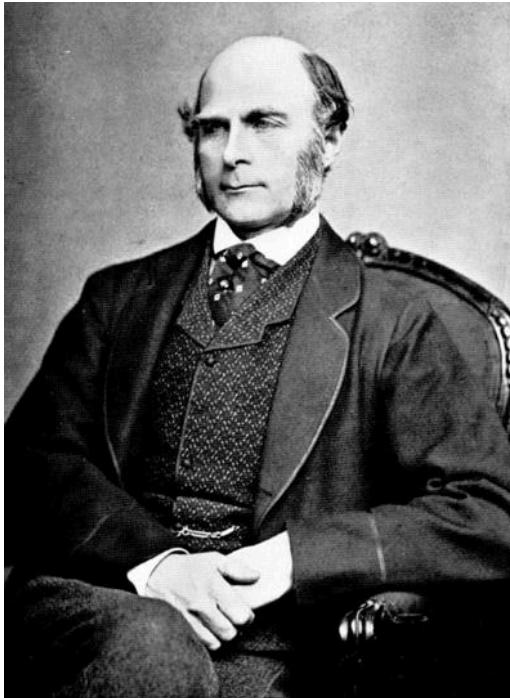
Humans only perceive the present

How do humans process temporal data?

Humans only perceive the present

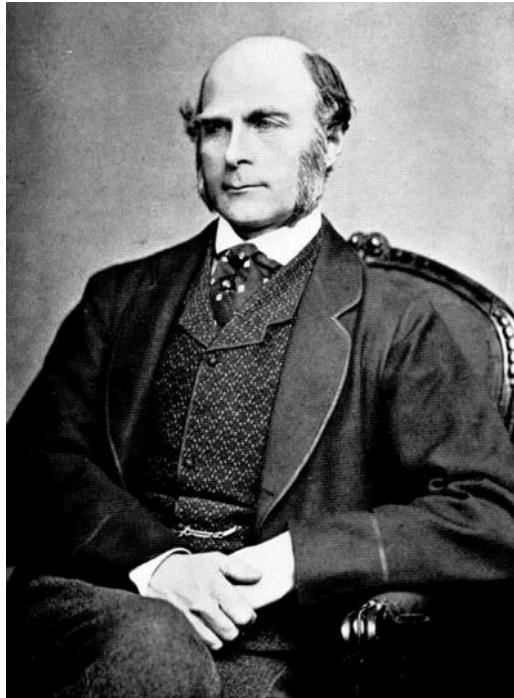
Humans process temporal data by storing and recalling memories

Francis Galton (1822-1911) photo composite memory

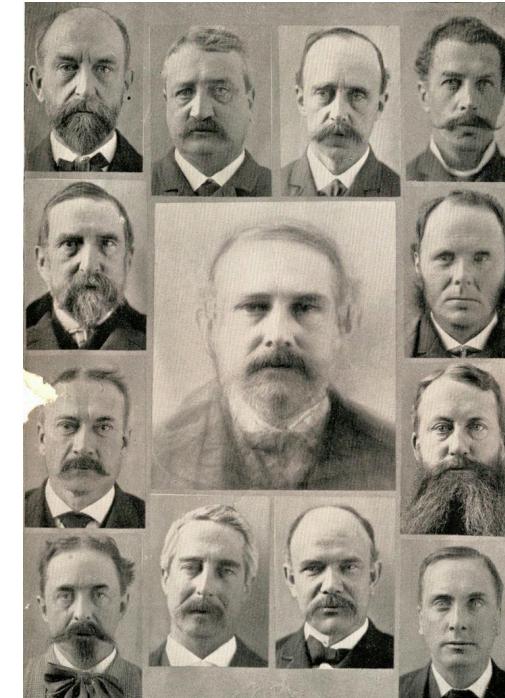


Review

Francis Galton (1822-1911)
photo composite memory



Composite photo of members of a
party



Task: Find a mathematical model of how our mind represents memories

Task: Find a mathematical model of how our mind represents memories

$X : \mathbb{R}^{h \times w}$ People you see at the party

Task: Find a mathematical model of how our mind represents memories

$$X : \mathbb{R}^{h \times w} \quad \text{People you see at the party}$$

$$H : \mathbb{R}^{h \times w} \quad \text{The image in your mind}$$

$$f : X^T \times \Theta \mapsto H$$

Task: Find a mathematical model of how our mind represents memories

$X : \mathbb{R}^{h \times w}$ People you see at the party

$H : \mathbb{R}^{h \times w}$ The image in your mind

$f : X^T \times \Theta \mapsto H$

Composite photography/memory uses a weighted sum

Task: Find a mathematical model of how our mind represents memories

$$X : \mathbb{R}^{h \times w} \quad \text{People you see at the party}$$

$$H : \mathbb{R}^{h \times w} \quad \text{The image in your mind}$$

$$f : X^T \times \Theta \mapsto H$$

Composite photography/memory uses a weighted sum

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \boldsymbol{\theta}^\top \overline{\mathbf{x}}_i$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \boldsymbol{\theta}^\top \overline{\mathbf{x}}_i$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \boldsymbol{\theta}^\top \overline{\mathbf{x}}_i$$

What if we see a new face?

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i$$

What if we see a new face?

$$f(\mathbf{x}, \boldsymbol{\theta}) = \left(\sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i \right) + \boldsymbol{\theta}^\top \bar{\mathbf{x}}_{\text{new}}$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i$$

What if we see a new face?

$$f(\mathbf{x}, \boldsymbol{\theta}) = \left(\sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i \right) + \boldsymbol{\theta}^\top \bar{\mathbf{x}}_{\text{new}}$$

We repeat the same process for each new face

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i$$

What if we see a new face?

$$f(\mathbf{x}, \boldsymbol{\theta}) = \left(\sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i \right) + \boldsymbol{\theta}^\top \bar{\mathbf{x}}_{\text{new}}$$

We repeat the same process for each new face

We can rewrite f as a **recurrent function**

Let us rewrite composite memory as a recurrent function

Let us rewrite composite memory as a recurrent function

$$f(\mathbf{x}, \boldsymbol{\theta}) = \underbrace{\left(\sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i \right)}_{\mathbf{h}} + \boldsymbol{\theta}^\top \bar{\mathbf{x}}_{\text{new}}$$

Let us rewrite composite memory as a recurrent function

$$f(x, \theta) = \underbrace{\left(\sum_{i=1}^T \theta^\top \bar{x}_i \right)}_h + \theta^\top \bar{x}_{\text{new}}$$

$$f(h, x, \theta) = h + \theta^\top \bar{x}$$

Let us rewrite composite memory as a recurrent function

$$f(\mathbf{x}, \boldsymbol{\theta}) = \underbrace{\left(\sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i \right)}_{\mathbf{h}} + \boldsymbol{\theta}^\top \bar{\mathbf{x}}_{\text{new}}$$

$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \mathbf{h} + \boldsymbol{\theta}^\top \bar{\mathbf{x}}$$

$$\mathbf{x} \in \mathbb{R}^{d_x}, \quad \mathbf{h} \in \mathbb{R}^{d_h}$$

Review

$$x \in \mathbb{R}^{d_x}, \quad h \in \mathbb{R}^{d_h} \quad f(h, x, \theta) = h + \theta^\top \bar{x}$$

Review

$$\mathbf{x} \in \mathbb{R}^{d_x}, \quad \mathbf{h} \in \mathbb{R}^{d_h} \quad f(\mathbf{h}, \mathbf{x}, \theta) = \mathbf{h} + \theta^\top \bar{\mathbf{x}}$$

$$\textcolor{red}{\mathbf{h}_1} = f(\mathbf{0}, \mathbf{x}_1, \theta) = \mathbf{0} + \theta^\top \bar{\mathbf{x}}_1$$

Review

$$\mathbf{x} \in \mathbb{R}^{d_x}, \quad \mathbf{h} \in \mathbb{R}^{d_h} \quad f(\mathbf{h}, \mathbf{x}, \theta) = \mathbf{h} + \theta^\top \bar{\mathbf{x}}$$

$$\textcolor{red}{\mathbf{h}_1} = f(\mathbf{0}, \mathbf{x}_1, \theta) = \mathbf{0} + \theta^\top \bar{\mathbf{x}}_1$$

$$\textcolor{green}{\mathbf{h}_2} = f(\textcolor{red}{\mathbf{h}_1}, \mathbf{x}_2, \theta) = \mathbf{h}_1 + \theta^\top \bar{\mathbf{x}}_2$$

Review

$$\mathbf{x} \in \mathbb{R}^{d_x}, \quad \mathbf{h} \in \mathbb{R}^{d_h} \quad f(\mathbf{h}, \mathbf{x}, \theta) = \mathbf{h} + \theta^\top \bar{\mathbf{x}}$$

$$\textcolor{red}{\mathbf{h}_1} = f(\mathbf{0}, \mathbf{x}_1, \theta) = \mathbf{0} + \theta^\top \bar{\mathbf{x}}_1$$

$$\textcolor{green}{\mathbf{h}_2} = f(\textcolor{red}{\mathbf{h}_1}, \mathbf{x}_2, \theta) = \mathbf{h}_1 + \theta^\top \bar{\mathbf{x}}_2$$

$$\mathbf{h}_3 = f(\textcolor{green}{\mathbf{h}_2}, \mathbf{x}_3, \theta) = \mathbf{h}_2 + \theta^\top \bar{\mathbf{x}}_3$$

Review

$$\mathbf{x} \in \mathbb{R}^{d_x}, \quad \mathbf{h} \in \mathbb{R}^{d_h} \quad f(\mathbf{h}, \mathbf{x}, \theta) = \mathbf{h} + \theta^\top \bar{\mathbf{x}}$$

$$\textcolor{red}{\mathbf{h}_1} = f(\mathbf{0}, \mathbf{x}_1, \theta) = \mathbf{0} + \theta^\top \bar{\mathbf{x}}_1$$

$$\textcolor{green}{\mathbf{h}_2} = f(\textcolor{red}{\mathbf{h}_1}, \mathbf{x}_2, \theta) = \mathbf{h}_1 + \theta^\top \bar{\mathbf{x}}_2$$

$$\mathbf{h}_3 = f(\textcolor{green}{\mathbf{h}_2}, \mathbf{x}_3, \theta) = \mathbf{h}_2 + \theta^\top \bar{\mathbf{x}}_3$$

⋮

$$\mathbf{h}_T = f(\mathbf{h}_{T-1}, \mathbf{x}_T, \theta) = \mathbf{h}_{T-1} + \theta^\top \bar{\mathbf{x}}_T$$

Review

$$\mathbf{x} \in \mathbb{R}^{d_x}, \quad \mathbf{h} \in \mathbb{R}^{d_h} \quad f(\mathbf{h}, \mathbf{x}, \theta) = \mathbf{h} + \theta^\top \bar{\mathbf{x}}$$

$$\textcolor{red}{\mathbf{h}_1} = f(\mathbf{0}, \mathbf{x}_1, \theta) = \mathbf{0} + \theta^\top \bar{\mathbf{x}}_1$$

$$\textcolor{green}{\mathbf{h}_2} = f(\textcolor{red}{\mathbf{h}_1}, \mathbf{x}_2, \theta) = \mathbf{h}_1 + \theta^\top \bar{\mathbf{x}}_2$$

$$\mathbf{h}_3 = f(\textcolor{green}{\mathbf{h}_2}, \mathbf{x}_3, \theta) = \mathbf{h}_2 + \theta^\top \bar{\mathbf{x}}_3$$

⋮

$$\mathbf{h}_T = f(\mathbf{h}_{T-1}, \mathbf{x}_T, \theta) = \mathbf{h}_{T-1} + \theta^\top \bar{\mathbf{x}}_T$$

Right now, our model remembers everything

Right now, our model remembers everything

But h is a fixed size, what if T is very large?

Right now, our model remembers everything

But h is a fixed size, what if T is very large?

If we keep adding and adding x into h , we will run out of space

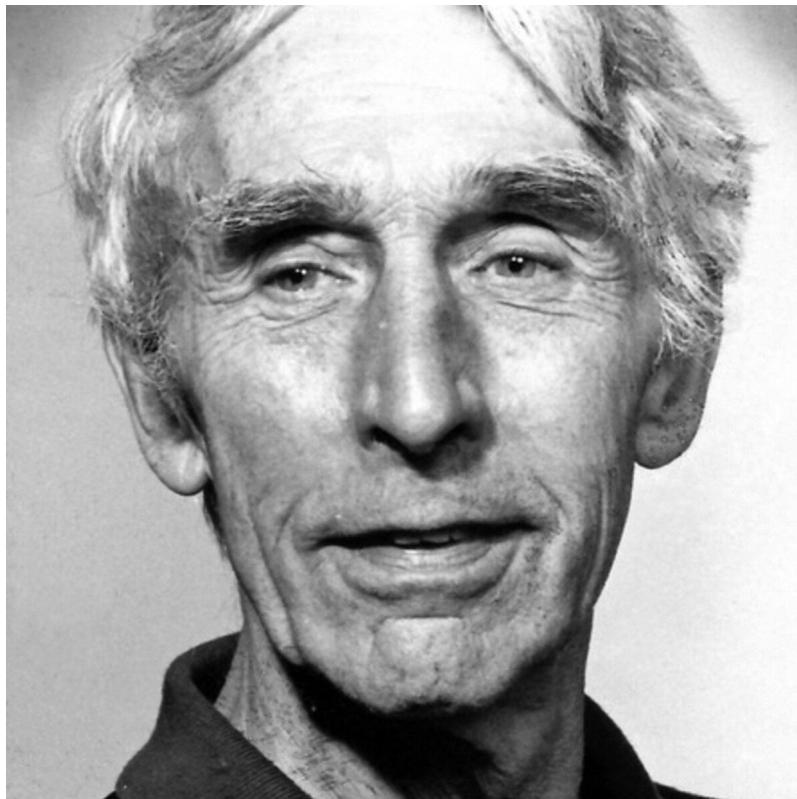
Right now, our model remembers everything

But h is a fixed size, what if T is very large?

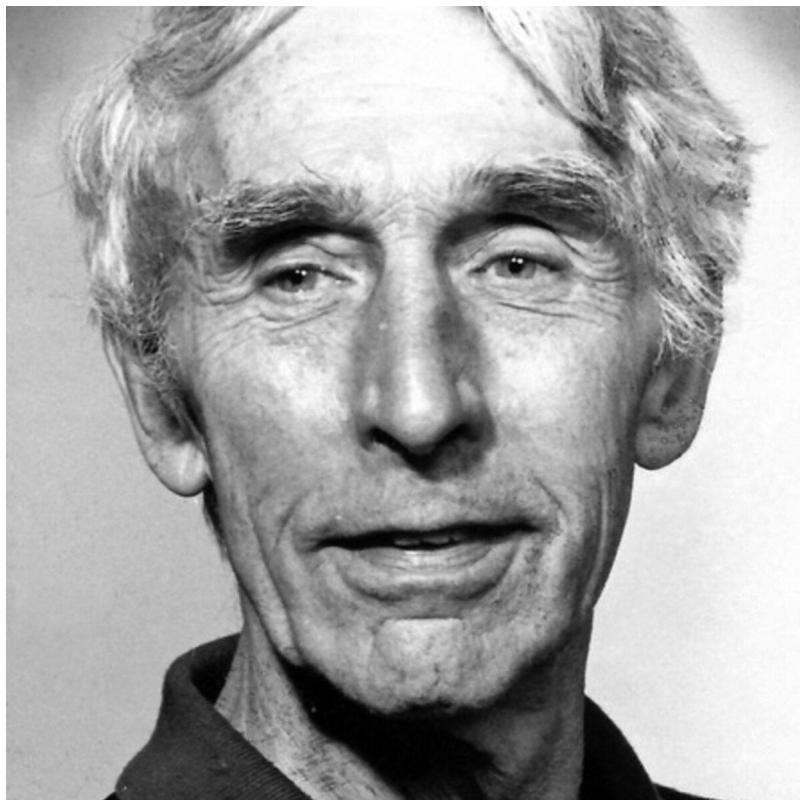
If we keep adding and adding x into h , we will run out of space

Humans forget old information

Murdock (1982)

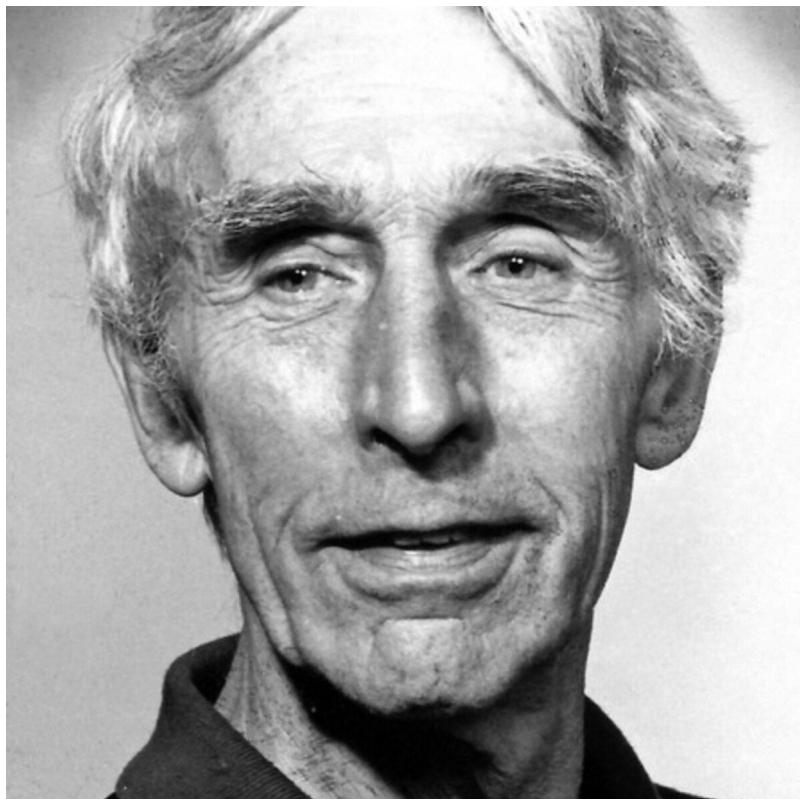


Murdock (1982)



$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \gamma \mathbf{h} + \boldsymbol{\theta}^\top \overline{\mathbf{x}}; \quad 0 < \gamma < 1$$

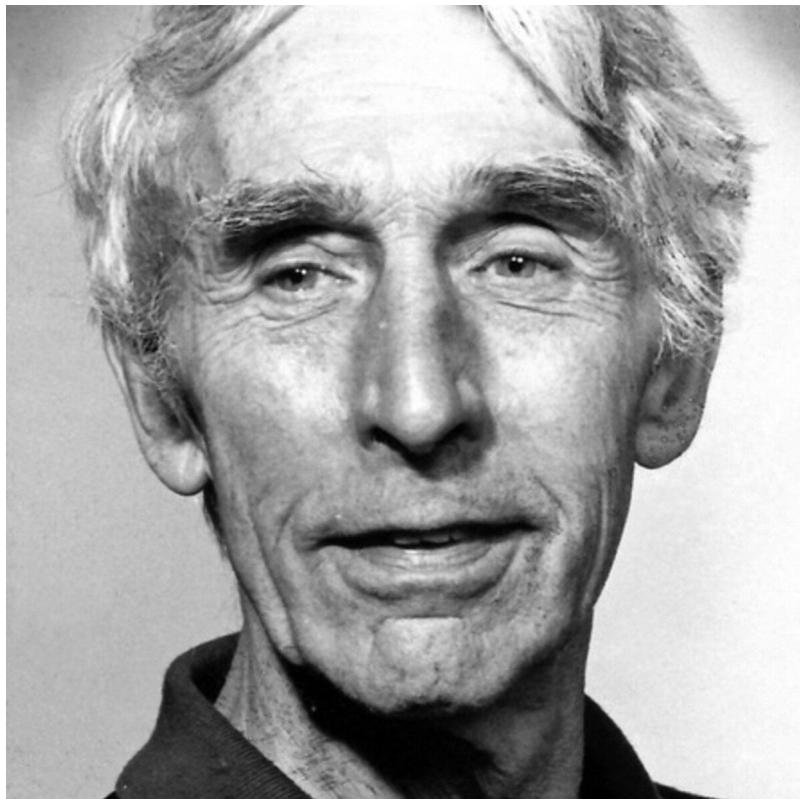
Murdock (1982)



$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \gamma \mathbf{h} + \boldsymbol{\theta}^\top \overline{\mathbf{x}}; \quad 0 < \gamma < 1$$

Key Idea: $\lim_{T \rightarrow \infty} \gamma^T = 0$

Murdock (1982)

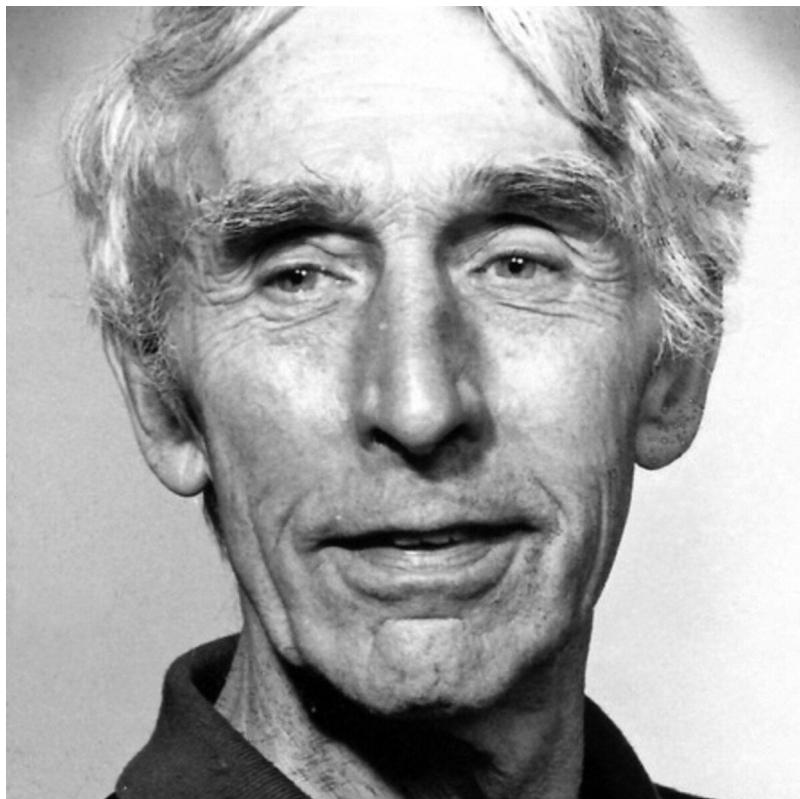


$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \gamma \mathbf{h} + \boldsymbol{\theta}^\top \overline{\mathbf{x}}; \quad 0 < \gamma < 1$$

Key Idea: $\lim_{T \rightarrow \infty} \gamma^T = 0$

Let $\gamma = 0.9$

Murdock (1982)



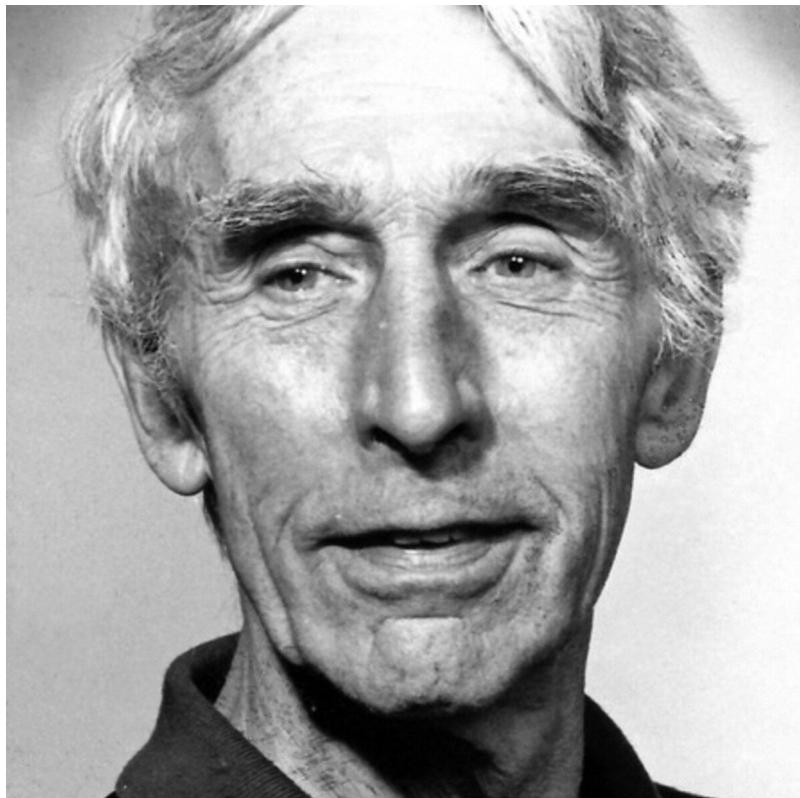
$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \gamma \mathbf{h} + \boldsymbol{\theta}^\top \bar{\mathbf{x}}; \quad 0 < \gamma < 1$$

Key Idea: $\lim_{T \rightarrow \infty} \gamma^T = 0$

Let $\gamma = 0.9$

$$0.9 \cdot 0.9 \cdot \mathbf{h} = 0.81\mathbf{h}$$

Murdock (1982)



$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \gamma \mathbf{h} + \boldsymbol{\theta}^\top \bar{\mathbf{x}}; \quad 0 < \gamma < 1$$

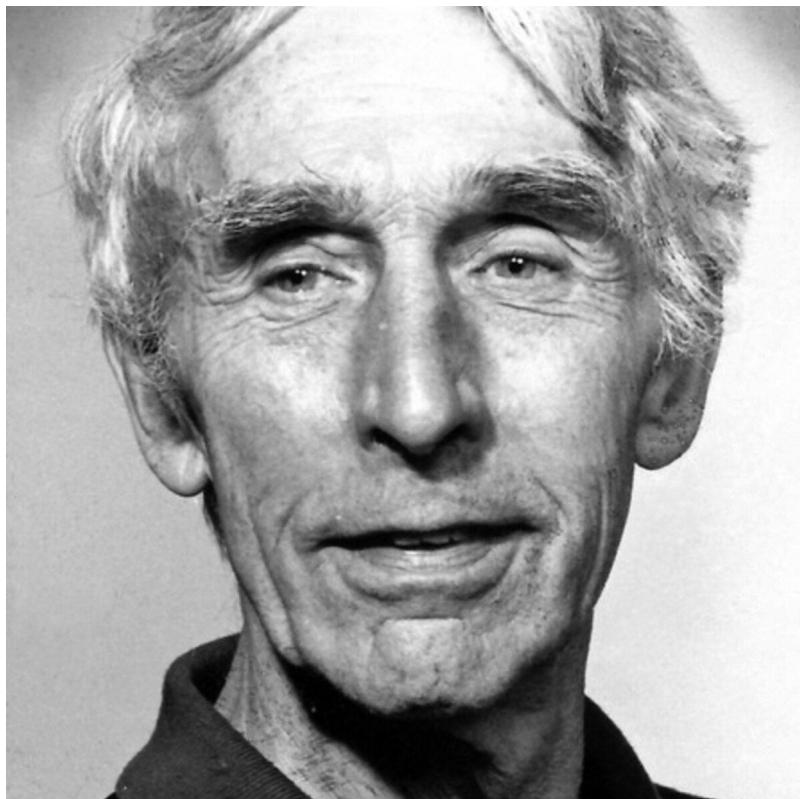
Key Idea: $\lim_{T \rightarrow \infty} \gamma^T = 0$

Let $\gamma = 0.9$

$$0.9 \cdot 0.9 \cdot \mathbf{h} = 0.81\mathbf{h}$$

$$0.9 \cdot 0.9 \cdot 0.9 \cdot \mathbf{h} = 0.72\mathbf{h}$$

Murdock (1982)



$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \gamma \mathbf{h} + \boldsymbol{\theta}^\top \bar{\mathbf{x}}; \quad 0 < \gamma < 1$$

Key Idea: $\lim_{T \rightarrow \infty} \gamma^T = 0$

Let $\gamma = 0.9$

$$0.9 \cdot 0.9 \cdot \mathbf{h} = 0.81\mathbf{h}$$

$$0.9 \cdot 0.9 \cdot 0.9 \cdot \mathbf{h} = 0.72\mathbf{h}$$

$$0.9 \cdot 0.9 \cdot 0.9 \cdot \dots \cdot \mathbf{h} = 0$$

$$h_T = \gamma^3 h_{T-3} + \gamma^2 \theta^\top \bar{x}_{T-2} + \gamma \theta^\top \bar{x}_{T-1} + \theta^\top \bar{x}_T$$

$$\boldsymbol{h}_T = \gamma^3 \boldsymbol{h}_{T-3} + \gamma^2 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_{T-2} + \gamma \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_{T-1} + \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_T$$

$$\boldsymbol{h}_T = \gamma^T \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_1 + \gamma^{T-1} \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_2 + \dots + \gamma^2 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_{T-2} + \gamma \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_{T-1} + \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_T$$

Our function f is just defined for a single X

Our function f is just defined for a single X

$$f : H \times X \times \Theta \mapsto H$$

Our function f is just defined for a single X

$$f : H \times X \times \Theta \mapsto H$$

To extend f to sequences, we scan f over the inputs

$$\text{scan}(f) : H \times X^T \times \Theta \mapsto H^T$$

Our function f is just defined for a single X

$$f : H \times X \times \Theta \mapsto H$$

To extend f to sequences, we scan f over the inputs

$$\text{scan}(f) : H \times X^T \times \Theta \mapsto H^T$$

$$\text{scan}(f)\left(h_0, \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_T \end{bmatrix}$$

$$f : H \times X \times \Theta \mapsto H, \quad \text{scan}(f) : H \times X^T \times \Theta \mapsto H^T$$

Review

$$f : H \times X \times \Theta \mapsto H, \quad \text{scan}(f) : H \times X^T \times \Theta \mapsto H^T$$

$$\text{scan}(f)\left(\mathbf{h}_0, \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta} \right) = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_T \end{bmatrix}$$

Review

$$f : H \times X \times \Theta \mapsto H, \quad \text{scan}(f) : H \times X^T \times \Theta \mapsto H^T$$

$$\text{scan}(f)\left(\mathbf{h}_0, \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \\ \vdots \\ \mathbf{h}_T \end{bmatrix}$$

$$\text{scan}(f)\left(\mathbf{h}_0, \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \begin{bmatrix} f(\mathbf{h}_0, \mathbf{x}_1, \boldsymbol{\theta}) \\ f(\mathbf{h}_1, \mathbf{x}_2, \boldsymbol{\theta}) \\ \vdots \\ f(\mathbf{h}_{T-1}, \mathbf{x}_T, \boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} f(\mathbf{h}_0, \mathbf{x}_1) \\ f(f(\mathbf{h}_0, \mathbf{x}_1), \mathbf{x}_2) \\ \vdots \\ f(\dots f(\mathbf{h}_0, \mathbf{x}_1) \dots, \mathbf{x}_T) \end{bmatrix}$$

Let g define our memory recall function

Let g define our memory recall function

$$g : H \times X \times \Theta \mapsto Y$$

Let g define our memory recall function

$$g : H \times X \times \Theta \mapsto Y$$

g searches your memories h using the input x , to produce output y

Let g define our memory recall function

$$g : H \times X \times \Theta \mapsto Y$$

g searches your memories h using the input x , to produce output y

x : “What is your favorite ice cream flavor?”

Let g define our memory recall function

$$g : H \times X \times \Theta \mapsto Y$$

g searches your memories h using the input x , to produce output y

x : “What is your favorite ice cream flavor?”

h : Everything you remember (hometown, birthday, etc)

Let g define our memory recall function

$$g : H \times X \times \Theta \mapsto Y$$

g searches your memories h using the input x , to produce output y

x : “What is your favorite ice cream flavor?”

h : Everything you remember (hometown, birthday, etc)

g : Searches your memories for ice cream information, and responds
“chocolate”

Step 1: Perform scan to find recurrent states

Step 1: Perform scan to find recurrent states

$$\begin{bmatrix} \mathbf{h}_1 \\ \vdots \\ \mathbf{h}_T \end{bmatrix} = \text{scan}(f) \left(\mathbf{h}_0, \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta_f \right)$$

Step 1: Perform scan to find recurrent states

$$\begin{bmatrix} \mathbf{h}_1 \\ \vdots \\ \mathbf{h}_T \end{bmatrix} = \text{scan}(f) \left(\mathbf{h}_0, \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta_f \right)$$

Step 2: Perform recall on recurrent states

Step 1: Perform scan to find recurrent states

$$\begin{bmatrix} \mathbf{h}_1 \\ \vdots \\ \mathbf{h}_T \end{bmatrix} = \text{scan}(f) \left(\mathbf{h}_0, \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta_f \right)$$

Step 2: Perform recall on recurrent states

$$\begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_T \end{bmatrix} = \begin{bmatrix} g(\mathbf{h}_1, \mathbf{x}_1, \theta_g) \\ \vdots \\ g(\mathbf{h}_T, \mathbf{x}_T, \theta_g) \end{bmatrix}$$

The simplest recurrent neural network is the **Elman Network**

The simplest recurrent neural network is the **Elman Network**

$$f(h, x, \theta) = \sigma(\theta_1^\top h + \theta_2^\top \bar{x})$$

The simplest recurrent neural network is the **Elman Network**

$$f(h, x, \theta) = \sigma(\theta_1^\top h + \theta_2^\top \bar{x})$$

$$g(h, x, \theta) = \theta_3^\top h$$

Add forgetting

$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^\top \mathbf{h} \odot f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) + \boldsymbol{\theta}_2^\top \bar{\mathbf{x}})$$

Add forgetting

$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^\top \mathbf{h} \odot f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) + \boldsymbol{\theta}_2^\top \bar{\mathbf{x}})$$

$$f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^\top \bar{\mathbf{x}} + \boldsymbol{\theta}_2^\top \mathbf{h})$$

Add forgetting

$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^\top \mathbf{h} \odot f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) + \boldsymbol{\theta}_2^\top \bar{\mathbf{x}})$$

$$f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^\top \bar{\mathbf{x}} + \boldsymbol{\theta}_2^\top \mathbf{h})$$

When $f_{\text{forget}} < 1$, we forget some of our memories!

Add forgetting

$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^\top \mathbf{h} \odot f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) + \boldsymbol{\theta}_2^\top \bar{\mathbf{x}})$$

$$f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^\top \bar{\mathbf{x}} + \boldsymbol{\theta}_2^\top \mathbf{h})$$

When $f_{\text{forget}} < 1$, we forget some of our memories!

Through gradient descent, neural network learns which memories to forget

Minimal gated unit (MGU) is a modern RNN

Minimal gated unit (MGU) is a modern RNN

MGU defines two helper functions

Minimal gated unit (MGU) is a modern RNN

MGU defines two helper functions

$$f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^\top \bar{\mathbf{x}} + \boldsymbol{\theta}_2^\top \mathbf{h})$$

Minimal gated unit (MGU) is a modern RNN

MGU defines two helper functions

$$f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^\top \bar{\mathbf{x}} + \boldsymbol{\theta}_2^\top \mathbf{h})$$

$$f_2(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_3^\top \bar{\mathbf{x}} + \boldsymbol{\theta}_4^\top f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) \odot \mathbf{h})$$

Minimal gated unit (MGU) is a modern RNN

MGU defines two helper functions

$$f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^\top \bar{\mathbf{x}} + \boldsymbol{\theta}_2^\top \mathbf{h})$$

$$f_2(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_3^\top \bar{\mathbf{x}} + \boldsymbol{\theta}_4^\top f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) \odot \mathbf{h})$$

$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) \odot \mathbf{h} + (1 - f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta})) \odot f_2(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta})$$

Minimal gated unit (MGU) is a modern RNN

MGU defines two helper functions

$$f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^\top \bar{\mathbf{x}} + \boldsymbol{\theta}_2^\top \mathbf{h})$$

$$f_2(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_3^\top \bar{\mathbf{x}} + \boldsymbol{\theta}_4^\top f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) \odot \mathbf{h})$$

$$f(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) = f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta}) \odot \mathbf{h} + (1 - f_{\text{forget}}(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta})) \odot f_2(\mathbf{h}, \mathbf{x}, \boldsymbol{\theta})$$

Left term forgets old, right term replaces forgotten memories

Review

Jax RNN https://colab.research.google.com/drive/147z7FNGyERV8oQ_4gZmxDVdeoNt0hKta#scrollTo=TUMonJ1u8Va

Agenda

1. Review
2. Compression
3. Autoencoders
4. Applications
5. Variational Modeling
6. VAE Implementation
7. Coding

Agenda

1. Review
2. **Compression**
3. Autoencoders
4. Applications
5. Variational Modeling
6. VAE Implementation
7. Coding

Compression



Compression



Question: You watch a film. How do you communicate information about the film with a friend?

Compression



Question: You watch a film. How do you communicate information about the film with a friend?

Answer: An ogre and donkey rescue a princess, discovering friendship and love along the way.

Compression



Question: You watch a film. How do you communicate information about the film with a friend?

Answer: An ogre and donkey rescue a princess, discovering friendship and love along the way.

Question: What is missing?

Compression



Question: You watch a film. How do you communicate information about the film with a friend?

Answer: An ogre and donkey rescue a princess, discovering friendship and love along the way.

Question: What is missing?

Answer: Shrek lives in a swamp, Lord Farquaad, dragons, etc

Compression

When you discuss films with friends, you summarize them

Compression

When you discuss films with friends, you summarize them

This is a form of **compression**

Compression

When you discuss films with friends, you summarize them

This is a form of **compression**

$$f \left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \right) = \text{Green ogre and donkey save princess}$$

Compression

When you discuss films with friends, you summarize them

This is a form of **compression**

$$f \left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \right) = \text{Green ogre and donkey save princess}$$

In compression, we reduce the size of data by removing information

Compression

When you discuss films with friends, you summarize them

This is a form of **compression**

$$f \left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \right) = \text{Green ogre and donkey save princess}$$

In compression, we reduce the size of data by removing information

Let us examine a more principled form of video compression

Compression

Shrek in 4k UHD:

$$X \in \mathbb{Z}_{255}^{3 \times 3840 \times 2160}, X^{90 \times 60 \times 24}$$

Compression

Shrek in 4k UHD:

$$X \in \mathbb{Z}_{255}^{3 \times 3840 \times 2160}, X^{90 \times 60 \times 24}$$

Question: How many GB?

Compression

Shrek in 4k UHD:

$$X \in \mathbb{Z}_{255}^{3 \times 3840 \times 2160}, X^{90 \times 60 \times 24}$$

Question: How many GB?

Answer: 3000 GB

Compression

Shrek in 4k UHD:

$$X \in \mathbb{Z}_{255}^{3 \times 3840 \times 2160}, X^{90 \times 60 \times 24}$$

Question: How many GB?

Answer: 3000 GB

But if you download films, you know they are smaller than 3 TB

Compression

Shrek in 4k UHD:

$$X \in \mathbb{Z}_{255}^{3 \times 3840 \times 2160}, X^{90 \times 60 \times 24}$$

Question: How many GB?

Answer: 3000 GB

But if you download films, you know they are smaller than 3 TB

Today, we use the H.264 video **encoder** to transform videos into a more compact representation

Compression

H.264 encoder selects 16×16 pixel blocks, estimates shift between frames, applies cosine transform, ...

Compression

H.264 encoder selects 16×16 pixel blocks, estimates shift between frames, applies cosine transform, ...

The result is an .mp4 file Z

Compression

H.264 encoder selects 16×16 pixel blocks, estimates shift between frames, applies cosine transform, ...

The result is an .mp4 file Z

$$f : X^t \mapsto Z \qquad Z \in \{0, 1\}^n$$

Compression

H.264 encoder selects 16×16 pixel blocks, estimates shift between frames, applies cosine transform, ...

The result is an .mp4 file Z

$$f : X^t \mapsto Z \qquad Z \in \{0, 1\}^n$$

Question: What is the size of Z in GB?

Compression

H.264 encoder selects 16×16 pixel blocks, estimates shift between frames, applies cosine transform, ...

The result is an .mp4 file Z

$$f : X^t \mapsto Z \qquad Z \in \{0, 1\}^n$$

Question: What is the size of Z in GB?

Answer: 60 GB, original size was 3000 GB

Compression

H.264 encoder selects 16×16 pixel blocks, estimates shift between frames, applies cosine transform, ...

The result is an .mp4 file Z

$$f : X^t \mapsto Z \qquad Z \in \{0, 1\}^n$$

Question: What is the size of Z in GB?

Answer: 60 GB, original size was 3000 GB

We achieve a compression ratio of $3000 \text{ GB} / 60 \text{ GB} = 50$

Compression

What happens on your computer when you watch Shrek?

Compression

What happens on your computer when you watch Shrek?

Download $z \in Z$ from the internet

Compression

What happens on your computer when you watch Shrek?

Download $z \in Z$ from the internet

$$Z \in \{0, 1\}^n$$

Compression

What happens on your computer when you watch Shrek?

Download $z \in Z$ from the internet $Z \in \{0, 1\}^n$

Information is no longer pixels, it is a string of bits

Compression

What happens on your computer when you watch Shrek?

Download $z \in Z$ from the internet $Z \in \{0, 1\}^n$

Information is no longer pixels, it is a string of bits

We must **decode** z back into pixels

Compression

What happens on your computer when you watch Shrek?

Download $z \in Z$ from the internet

$$Z \in \{0, 1\}^n$$

Information is no longer pixels, it is a string of bits

We must **decode** z back into pixels

We need to undo (invert) the encoder f

$$f : X^t \mapsto Z$$

$$f^{-1} : Z \mapsto X^t$$

Compression

What happens on your computer when you watch Shrek?

Download $z \in Z$ from the internet

$$Z \in \{0, 1\}^n$$

Information is no longer pixels, it is a string of bits

We must **decode** z back into pixels

We need to undo (invert) the encoder f

$$f : X^t \mapsto Z$$

$$f^{-1} : Z \mapsto X^t$$

Your CPU has a H.264 decoder built in to make this fast

Compression

Compression may be **lossy**

Compression

Compression may be **lossy** or **lossless**

Compression

Compression may be **lossy**

or **lossless**



Compression

Compression may be **lossy** or **lossless**



Question: Which is H.264?

Agenda

1. Review
2. **Compression**
3. Autoencoders
4. Applications
5. Variational Modeling
6. VAE Implementation
7. Coding

Agenda

1. Review
2. Compression
3. **Autoencoders**
4. Applications
5. Variational Modeling
6. VAE Implementation
7. Coding

Autoencoders

Encoders and decoders for images, videos, and music are functions

Autoencoders

Encoders and decoders for images, videos, and music are functions

Neural networks can represent any continuous function

Autoencoders

Encoders and decoders for images, videos, and music are functions

Neural networks can represent any continuous function

Let us learn neural network encoders and decoders

Autoencoders

Encoders and decoders for images, videos, and music are functions

Neural networks can represent any continuous function

Let us learn neural network encoders and decoders

$$f : X \times \Theta \mapsto Z$$

Autoencoders

Encoders and decoders for images, videos, and music are functions

Neural networks can represent any continuous function

Let us learn neural network encoders and decoders

$$f : X \times \Theta \mapsto Z$$

$$f^{-1} : Z \times \Theta \mapsto X$$

Autoencoders

Encoders and decoders for images, videos, and music are functions

Neural networks can represent any continuous function

Let us learn neural network encoders and decoders

$$f : X \times \Theta \mapsto Z$$

$$f^{-1} : Z \times \Theta \mapsto X$$

We call this an **autoencoder**

Autoencoders

Encoders and decoders for images, videos, and music are functions

Neural networks can represent any continuous function

Let us learn neural network encoders and decoders

$$f : X \times \Theta \mapsto Z$$

$$f^{-1} : Z \times \Theta \mapsto X$$

We call this an **autoencoder**

Notice there is no Y this time

Autoencoders

Encoders and decoders for images, videos, and music are functions

Neural networks can represent any continuous function

Let us learn neural network encoders and decoders

$$f : X \times \Theta \mapsto Z$$

$$f^{-1} : Z \times \Theta \mapsto X$$

We call this an **autoencoder**

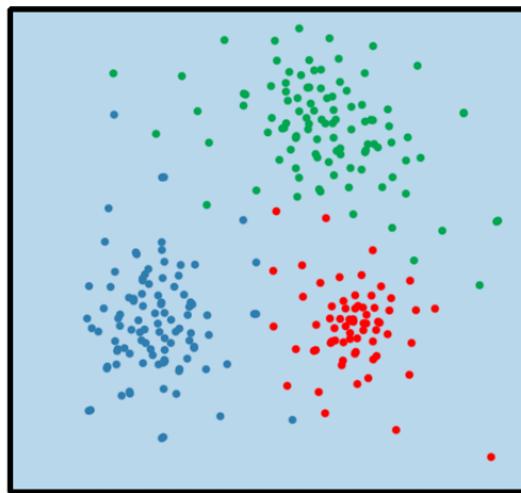
Notice there is no Y this time

Training autoencoders is different than what we have seen before

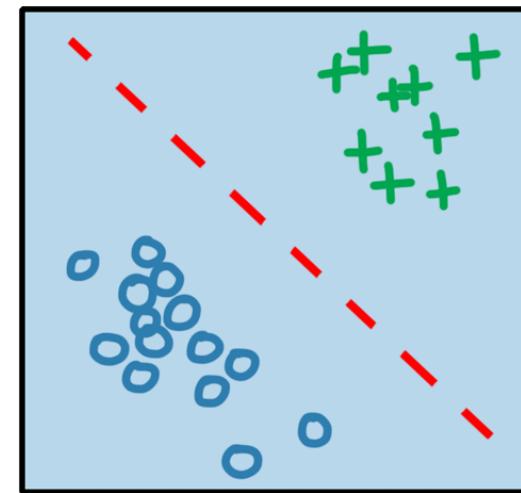
Autoencoders

machine learning

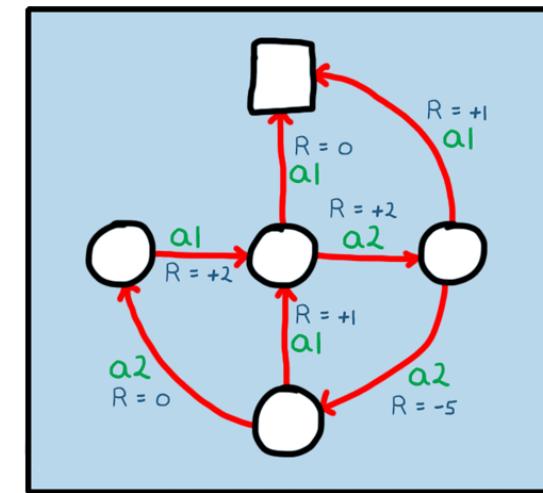
unsupervised
learning



supervised
learning



reinforcement
learning



Autoencoders

In supervised learning, humans provide the model with **inputs X** and corresponding **outputs Y**

Autoencoders

In supervised learning, humans provide the model with **inputs X** and corresponding **outputs Y**

$$\mathbf{X} = [x_{[1]} \ x_{[2]} \ \dots \ x_{[n]}]^\top \quad \mathbf{Y} = [y_{[1]} \ y_{[2]} \ \dots \ y_{[n]}]^\top$$

Autoencoders

In supervised learning, humans provide the model with **inputs X** and corresponding **outputs Y**

$$\mathbf{X} = [x_{[1]} \ x_{[2]} \ \dots \ x_{[n]}]^\top \quad \mathbf{Y} = [y_{[1]} \ y_{[2]} \ \dots \ y_{[n]}]^\top$$

In unsupervised learning, humans only provide **input**

Autoencoders

In supervised learning, humans provide the model with **inputs X** and corresponding **outputs Y**

$$\mathbf{X} = [x_{[1]} \ x_{[2]} \ \dots \ x_{[n]}]^\top \quad \mathbf{Y} = [y_{[1]} \ y_{[2]} \ \dots \ y_{[n]}]^\top$$

In unsupervised learning, humans only provide **input**

$$\mathbf{X} = [x_{[1]} \ x_{[2]} \ \dots \ x_{[n]}]^\top$$

Autoencoders

In supervised learning, humans provide the model with **inputs X** and corresponding **outputs Y**

$$\mathbf{X} = [x_{[1]} \ x_{[2]} \ \dots \ x_{[n]}]^\top \quad \mathbf{Y} = [y_{[1]} \ y_{[2]} \ \dots \ y_{[n]}]^\top$$

In unsupervised learning, humans only provide **input**

$$\mathbf{X} = [x_{[1]} \ x_{[2]} \ \dots \ x_{[n]}]^\top$$

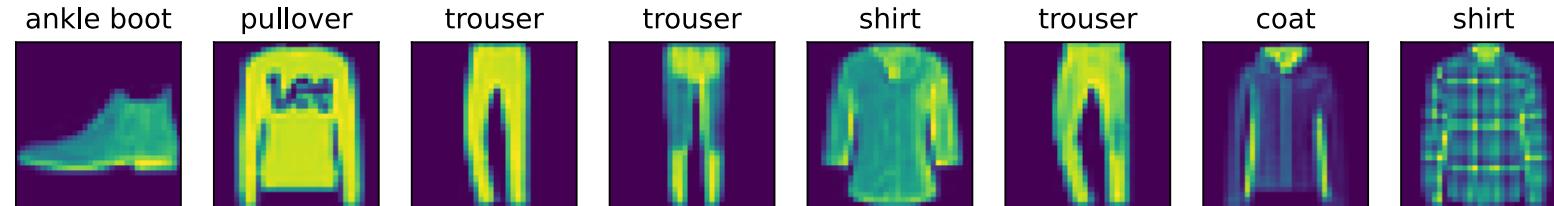
The training algorithm may generate labels

Autoencoders

Task: Compress images for your clothing website to save on costs

Autoencoders

Task: Compress images for your clothing website to save on costs



$$d_x : 28 \times 28$$

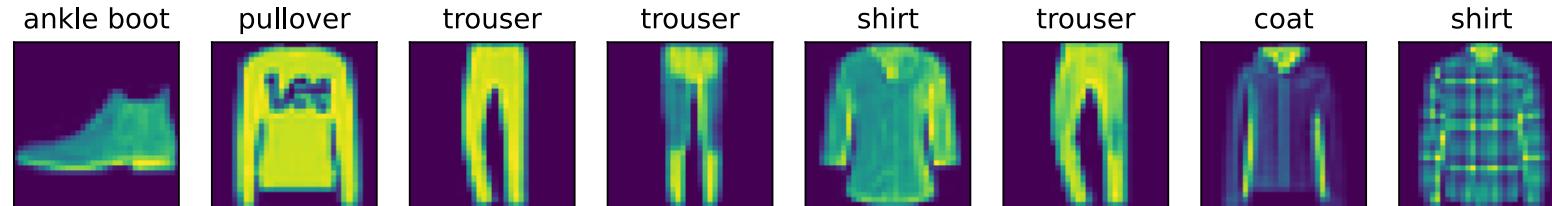
$$X \in [0, 1]^{d_x}$$

$$d_z : 4$$

$$Z \in \mathbb{R}^{d_z}$$

Autoencoders

Task: Compress images for your clothing website to save on costs



$$d_x : 28 \times 28$$

$$d_z : 4$$

$$X \in [0, 1]^{d_x}$$

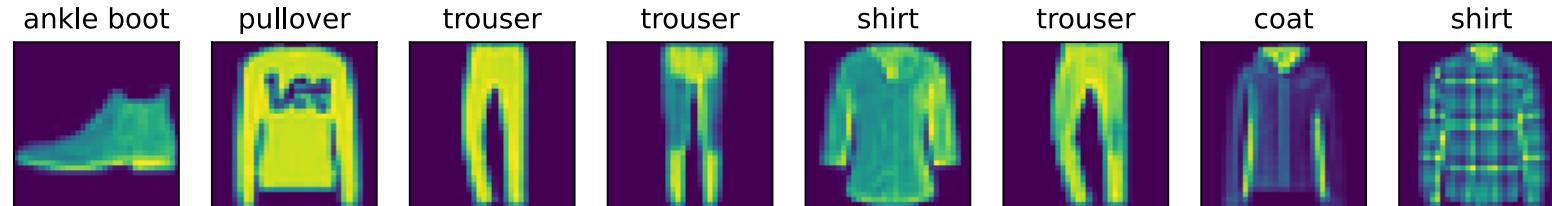
$$Z \in \mathbb{R}^{d_z}$$

$$f : X \times \Theta \mapsto Z$$

$$f^{-1} : Z \times \Theta \mapsto X$$

Autoencoders

Task: Compress images for your clothing website to save on costs



$$d_x : 28 \times 28$$

$$d_z : 4$$

$$X \in [0, 1]^{d_x}$$

$$Z \in \mathbb{R}^{d_z}$$

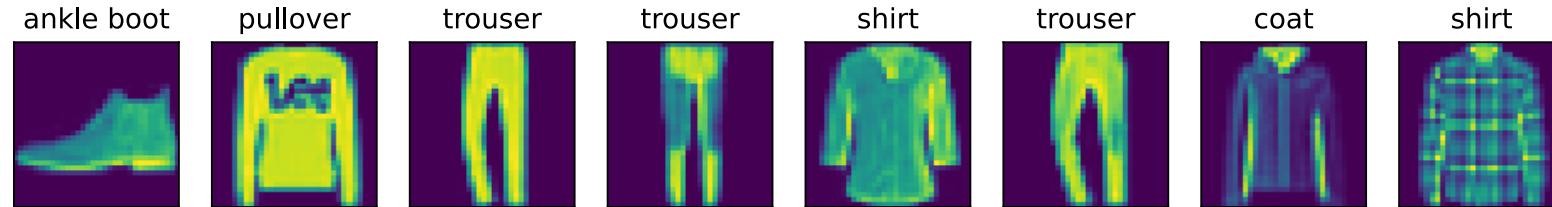
$$f : X \times \Theta \mapsto Z$$

$$f^{-1} : Z \times \Theta \mapsto X$$

What is the structure of f, f^{-1} ?

Autoencoders

Task: Compress images for your clothing website to save on costs



$$d_x : 28 \times 28$$

$$X \in [0, 1]^{d_x}$$

$$f : X \times \Theta \mapsto Z$$

What is the structure of f, f^{-1} ?

$$d_z : 4$$

$$Z \in \mathbb{R}^{d_z}$$

$$f^{-1} : Z \times \Theta \mapsto X$$

How do we find θ ?

Autoencoders

Let us find f , then find the inverse f^{-1}

Autoencoders

Let us find f , then find the inverse f^{-1}

Start with a perceptron

Autoencoders

Let us find f , then find the inverse f^{-1}

Start with a perceptron

$$f(x, \theta) = \sigma(\theta^\top \bar{x}); \quad \theta \in \mathbb{R}^{d_x \times d_z}$$

Autoencoders

Let us find f , then find the inverse f^{-1}

Start with a perceptron

$$f(x, \theta) = \sigma(\theta^\top \bar{x}); \quad \theta \in \mathbb{R}^{d_x \times d_z}$$

$$z = \sigma(\theta^\top \bar{x})$$

Autoencoders

Let us find f , then find the inverse f^{-1}

Start with a perceptron

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \bar{\boldsymbol{x}}); \quad \boldsymbol{\theta} \in \mathbb{R}^{d_x \times d_z}$$

$$\boldsymbol{z} = \sigma(\boldsymbol{\theta}^\top \bar{\boldsymbol{x}})$$

Solve for \boldsymbol{x} to find the inverse

Autoencoders

Let us find f , then find the inverse f^{-1}

Start with a perceptron

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \bar{\boldsymbol{x}}); \quad \boldsymbol{\theta} \in \mathbb{R}^{d_x \times d_z}$$

$$\boldsymbol{z} = \sigma(\boldsymbol{\theta}^\top \bar{\boldsymbol{x}})$$

Solve for \boldsymbol{x} to find the inverse

$$\sigma^{-1}(\boldsymbol{z}) = \sigma^{-1}(\sigma(\boldsymbol{\theta}^\top \bar{\boldsymbol{x}}))$$

Autoencoders

Let us find f , then find the inverse f^{-1}

Start with a perceptron

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \bar{\boldsymbol{x}}); \quad \boldsymbol{\theta} \in \mathbb{R}^{d_x \times d_z}$$

$$\boldsymbol{z} = \sigma(\boldsymbol{\theta}^\top \bar{\boldsymbol{x}})$$

Solve for \boldsymbol{x} to find the inverse

$$\sigma^{-1}(\boldsymbol{z}) = \sigma^{-1}(\sigma(\boldsymbol{\theta}^\top \bar{\boldsymbol{x}}))$$

$$\sigma^{-1}(\boldsymbol{z}) = \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}$$

Autoencoders

$$\sigma^{-1}(z) = \theta^\top \bar{x}$$

Autoencoders

$$\sigma^{-1}(z) = \theta^\top \bar{x}$$

$$(\theta^\top)^{-1} \sigma^{-1}(z) = (\theta^\top)^{-1} \theta^\top \bar{x}$$

Autoencoders

$$\sigma^{-1}(z) = \theta^\top \bar{x}$$

$$(\theta^\top)^{-1} \sigma^{-1}(z) = (\theta^\top)^{-1} \theta^\top \bar{x}$$

$$(\theta^\top)^{-1} \sigma^{-1}(z) = I \bar{x}$$

Autoencoders

$$\sigma^{-1}(z) = \theta^\top \bar{x}$$

$$(\theta^\top)^{-1} \sigma^{-1}(z) = (\theta^\top)^{-1} \theta^\top \bar{x}$$

$$(\theta^\top)^{-1} \sigma^{-1}(z) = I \bar{x}$$

$$(\theta^\top)^{-1} \sigma^{-1}(z) = \bar{x}$$

Autoencoders

$$\sigma^{-1}(z) = \theta^\top \bar{x}$$

$$(\theta^\top)^{-1} \sigma^{-1}(z) = (\theta^\top)^{-1} \theta^\top \bar{x}$$

$$(\theta^\top)^{-1} \sigma^{-1}(z) = I \bar{x}$$

$$(\theta^\top)^{-1} \sigma^{-1}(z) = \bar{x}$$

$$f^{-1}(z, \theta) = (\theta^\top)^{-1} \sigma^{-1}(z)$$

Autoencoders

$$f(x, \theta) = \sigma(\theta^\top \bar{x}) \quad f^{-1}(z, \theta) = (\theta^\top)^{-1} \sigma^{-1}(z) \quad \theta \in \mathbb{R}^{d_x \times d_z}$$

Autoencoders

$$f(x, \theta) = \sigma(\theta^\top \bar{x}) \quad f^{-1}(z, \theta) = (\theta^\top)^{-1} \sigma^{-1}(z) \quad \theta \in \mathbb{R}^{d_x \times d_z}$$

Question: Any issues?

Hint: What if $d_x \neq d_z$?

Autoencoders

$$f(x, \theta) = \sigma(\theta^\top \bar{x}) \quad f^{-1}(z, \theta) = (\theta^\top)^{-1} \sigma^{-1}(z) \quad \theta \in \mathbb{R}^{d_x \times d_z}$$

Question: Any issues?

Hint: What if $d_x \neq d_z$?

Answer: Can only invert square matrices, θ^\top only invertible if $d_z = d_x$

Autoencoders

$$f(x, \theta) = \sigma(\theta^\top \bar{x}) \quad f^{-1}(z, \theta) = (\theta^\top)^{-1} \sigma^{-1}(z) \quad \theta \in \mathbb{R}^{d_x \times d_z}$$

Question: Any issues?

Hint: What if $d_x \neq d_z$?

Answer: Can only invert square matrices, θ^\top only invertible if $d_z = d_x$

Question: What kind of compression can we achieve if $d_z = d_x$?

Autoencoders

$$f(x, \theta) = \sigma(\theta^\top \bar{x}) \quad f^{-1}(z, \theta) = (\theta^\top)^{-1} \sigma^{-1}(z) \quad \theta \in \mathbb{R}^{d_x \times d_z}$$

Question: Any issues?

Hint: What if $d_x \neq d_z$?

Answer: Can only invert square matrices, θ^\top only invertible if $d_z = d_x$

Question: What kind of compression can we achieve if $d_z = d_x$?

Answer: None! We need $d_z < d_x$ for compression

Autoencoders

$$f(x, \theta) = \sigma(\theta^\top \bar{x}) \quad f^{-1}(z, \theta) = (\theta^\top)^{-1} \sigma^{-1}(z) \quad \theta \in \mathbb{R}^{d_x \times d_z}$$

Question: Any issues?

Hint: What if $d_x \neq d_z$?

Answer: Can only invert square matrices, θ^\top only invertible if $d_z = d_x$

Question: What kind of compression can we achieve if $d_z = d_x$?

Answer: None! We need $d_z < d_x$ for compression

Look for another solution

Autoencoders

Let us try another way

Autoencoders

Let us try another way

$$z = f(x, \theta_e) = \sigma(\theta_e^\top \bar{x})$$

$$x = f^{-1}(z, \theta_d) = \sigma(\theta_d^\top \bar{z})$$

Autoencoders

Let us try another way

$$z = f(x, \theta_e) = \sigma(\theta_e^\top \bar{x})$$

$$x = f^{-1}(z, \theta_d) = \sigma(\theta_d^\top \bar{z})$$

What if we plug z into the second equation?

Autoencoders

Let us try another way

$$z = f(x, \theta_e) = \sigma(\theta_e^\top \bar{x})$$

$$x = f^{-1}(z, \theta_d) = \sigma(\theta_d^\top \bar{z})$$

What if we plug z into the second equation?

$$x = f^{-1}(f(x, \theta_e), \theta_d) = \sigma(\theta_d^\top \sigma(\theta_e^\top \bar{x}))$$

Autoencoders

$$\boldsymbol{x} = f^{-1}(f(\boldsymbol{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d) = \sigma(\boldsymbol{\theta}_d^\top \boldsymbol{\sigma}(\boldsymbol{\theta}_e^\top \boldsymbol{\bar{x}}))$$

Autoencoders

$$\boldsymbol{x} = f^{-1}(f(\boldsymbol{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d) = \sigma(\boldsymbol{\theta}_d^\top \boldsymbol{\sigma}(\boldsymbol{\theta}_e^\top \boldsymbol{x}))$$

More generally, f, f^{-1} may be any neural network

Autoencoders

$$\boldsymbol{x} = f^{-1}(f(\boldsymbol{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d) = \sigma(\boldsymbol{\theta}_d^\top \boldsymbol{\sigma}(\boldsymbol{\theta}_e^\top \bar{\boldsymbol{x}}))$$

More generally, f, f^{-1} may be any neural network

$$\boldsymbol{x} = f^{-1}(f(\boldsymbol{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)$$

Autoencoders

$$\boldsymbol{x} = f^{-1}(f(\boldsymbol{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d) = \sigma(\boldsymbol{\theta}_d^\top \boldsymbol{\sigma}(\boldsymbol{\theta}_e^\top \boldsymbol{x}))$$

More generally, f, f^{-1} may be any neural network

$$\boldsymbol{x} = f^{-1}(f(\boldsymbol{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)$$

Turn this into a loss function using the square error

Autoencoders

$$\mathbf{x} = f^{-1}(f(\mathbf{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d) = \sigma(\boldsymbol{\theta}_d^\top \boldsymbol{\sigma}(\boldsymbol{\theta}_e^\top \mathbf{x}))$$

More generally, f, f^{-1} may be any neural network

$$\mathbf{x} = f^{-1}(f(\mathbf{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)$$

Turn this into a loss function using the square error

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^{d_x} \left(x_j - f^{-1}(f(\mathbf{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)_j \right)^2$$

Autoencoders

$$\boldsymbol{x} = f^{-1}(f(\boldsymbol{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d) = \sigma(\boldsymbol{\theta}_d^\top \boldsymbol{\sigma}(\boldsymbol{\theta}_e^\top \boldsymbol{x}))$$

More generally, f, f^{-1} may be any neural network

$$\boldsymbol{x} = f^{-1}(f(\boldsymbol{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)$$

Turn this into a loss function using the square error

$$\mathcal{L}(\boldsymbol{x}, \boldsymbol{\theta}) = \sum_{j=1}^{d_x} \left(x_j - f^{-1}(f(\boldsymbol{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)_j \right)^2$$

Forces the networks to compress and reconstruct \boldsymbol{x}

Autoencoders

$$\mathcal{L}(x, \theta) = \sum_{j=1}^{d_x} \left(x_j - f^{-1}(f(x, \theta_e), \theta_d)_j \right)^2$$

Autoencoders

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^{d_x} \left(x_j - f^{-1}(f(\mathbf{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)_j \right)^2$$

Define over the entire dataset

$$\mathcal{L}(\mathbf{X}, \boldsymbol{\theta}) = \sum_{i=1}^n \sum_{j=1}^{d_x} \left(x_{[i],j} - f^{-1}(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)_j \right)^2$$

Autoencoders

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^{d_x} \left(x_j - f^{-1}(f(\mathbf{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)_j \right)^2$$

Define over the entire dataset

$$\mathcal{L}(\mathbf{X}, \boldsymbol{\theta}) = \sum_{i=1}^n \sum_{j=1}^{d_x} \left(x_{[i],j} - f^{-1}(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)_j \right)^2$$

We call this the **reconstruction loss**

Autoencoders

$$\mathcal{L}(x, \theta) = \sum_{j=1}^{d_x} \left(x_j - f^{-1}(f(x, \theta_e), \theta_d)_j \right)^2$$

Define over the entire dataset

$$\mathcal{L}(X, \theta) = \sum_{i=1}^n \sum_{j=1}^{d_x} \left(x_{[i],j} - f^{-1}(f(x_{[i]}, \theta_e), \theta_d)_j \right)^2$$

We call this the **reconstruction loss**

It is an unsupervised loss because we only provide X and not Y !

Autoencoders

First coding exercise

https://colab.research.google.com/drive/1UyR_W6NDIujaJXYlHZh6O3NfaCAMscpH#scrollTo=nmyQ8aE2pSbb

<https://www.youtube.com/watch?v=UZDiGooFs54>

Agenda

1. Review
2. Compression
3. **Autoencoders**
4. Applications
5. Variational Modeling
6. VAE Implementation
7. Coding

Agenda

1. Review
2. Compression
3. Autoencoders
4. **Applications**
5. Variational Modeling
6. VAE Implementation
7. Coding

Applications

We can use autoencoders for more than compression

Applications

We can use autoencoders for more than compression

We can make **denoising autoencoders** that remove noise

Applications

We can use autoencoders for more than compression

We can make **denoising autoencoders** that remove noise



Applications

Generate some noise

$$\varepsilon \sim \mathcal{N}(\mu, \Sigma)$$

Applications

Generate some noise

$$\varepsilon \sim \mathcal{N}(\mu, \Sigma)$$

Add noise to the image

$$x + \varepsilon$$

Applications

Generate some noise

$$\varepsilon \sim \mathcal{N}(\mu, \Sigma)$$

Add noise to the image

$$x + \varepsilon$$

Original loss $\mathcal{L}(X, \theta) = \sum_{i=1}^n \sum_{j=1}^{d_x} \left(x_{[i],j} - f^{-1} \left(f(x_{[i]}, \theta_e), \theta_d \right)_j \right)^2$

Applications

Generate some noise

$$\varepsilon \sim \mathcal{N}(\mu, \Sigma)$$

Add noise to the image

$$x + \varepsilon$$

Original loss $\mathcal{L}(X, \theta) = \sum_{i=1}^n \sum_{j=1}^{d_x} \left(x_{[i],j} - f^{-1}\left(f(x_{[i]}, \theta_e), \theta_d\right)_j \right)^2$

Denoising loss $\mathcal{L}(X, \theta) = \sum_{i=1}^n \sum_{j=1}^{d_x} \left(x_{[i],j} - f^{-1}\left(f(x_{[i]} + \varepsilon, \theta_e), \theta_d\right)_j \right)^2$

Applications

Generate some noise

$$\varepsilon \sim \mathcal{N}(\mu, \Sigma)$$

Add noise to the image

$$x + \varepsilon$$

Original loss $\mathcal{L}(X, \theta) = \sum_{i=1}^n \sum_{j=1}^{d_x} \left(x_{[i],j} - f^{-1}\left(f(x_{[i]}, \theta_e), \theta_d\right)_j \right)^2$

Denoising loss $\mathcal{L}(X, \theta) = \sum_{i=1}^n \sum_{j=1}^{d_x} \left(x_{[i],j} - f^{-1}\left(f(x_{[i]} + \varepsilon, \theta_e), \theta_d\right)_j \right)^2$

Autoencoder will learn to remove noise when reconstructing image

Applications

We can add camera blur too

Applications

We can add camera blur too



Applications

$$\text{blur}(x + \epsilon)$$

Applications

$$\text{blur}(x + \varepsilon)$$

Denoising deblur loss

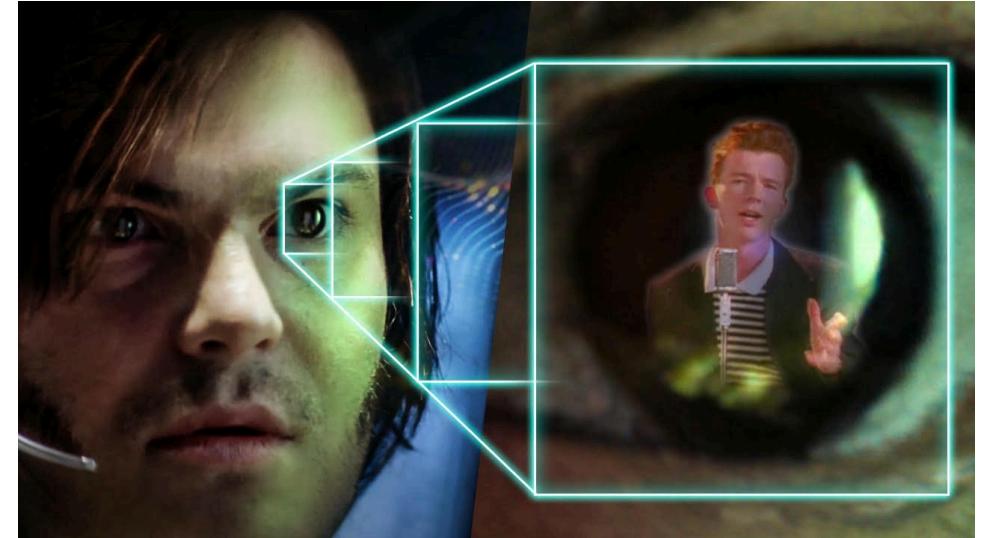
$$\mathcal{L}(\mathbf{X}, \boldsymbol{\theta}) = \sum_{i=1}^n \sum_{j=1}^{d_x} \left(x_{[i],j} - f^{-1} \left(f \left(\text{blur} \left(\mathbf{x}_{[i]} + \varepsilon \right), \boldsymbol{\theta}_e \right), \boldsymbol{\theta}_d \right)_j \right)^2$$

Applications

Now we can “enhance” images like in crime tv shows

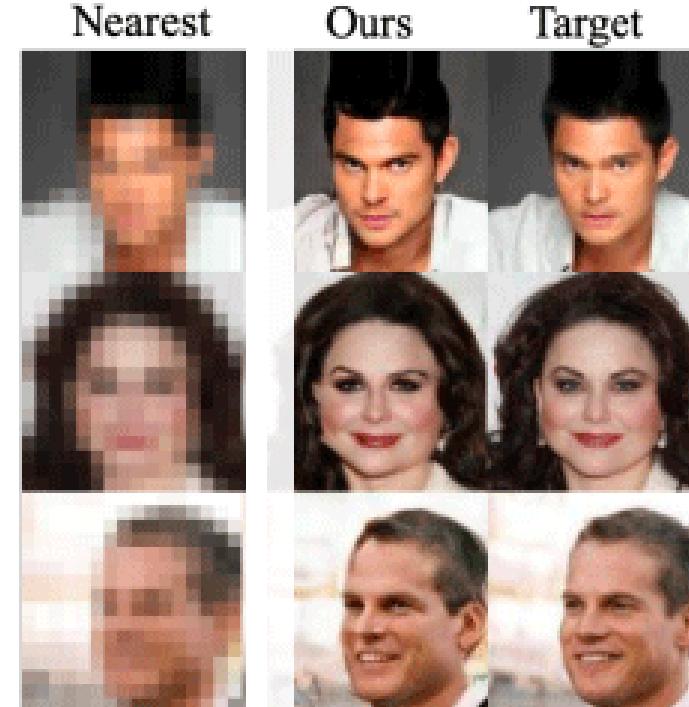
Applications

Now we can “enhance” images like in crime tv shows



Applications

We can deblur faces from security cameras



Applications

We can even hide parts of the image

Applications

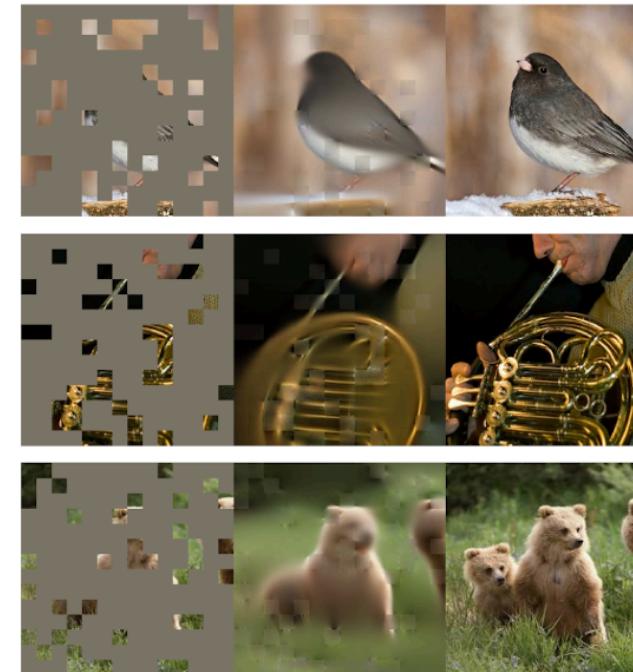
We can even hide parts of the image

A **masked autoencoder** will reconstruct the missing data

Applications

We can even hide parts of the image

A **masked autoencoder** will reconstruct the missing data



Applications

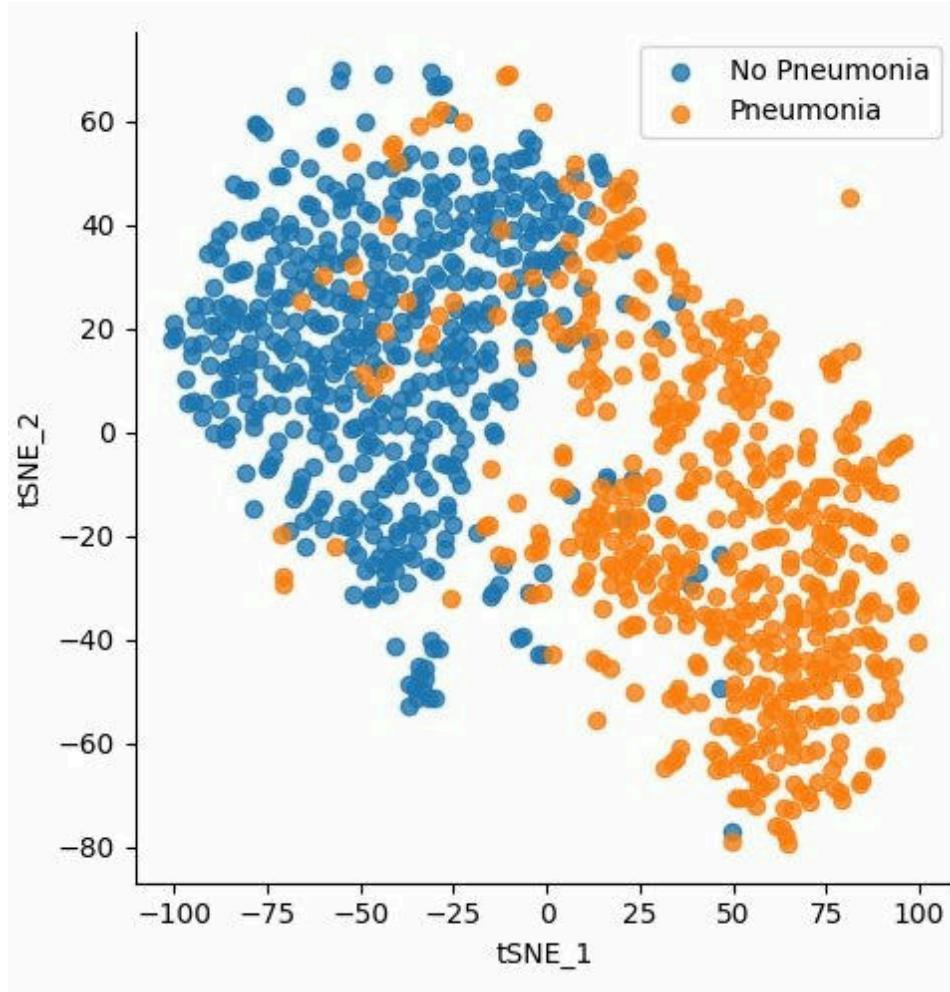
What is happening here? How can these models do this?

Applications

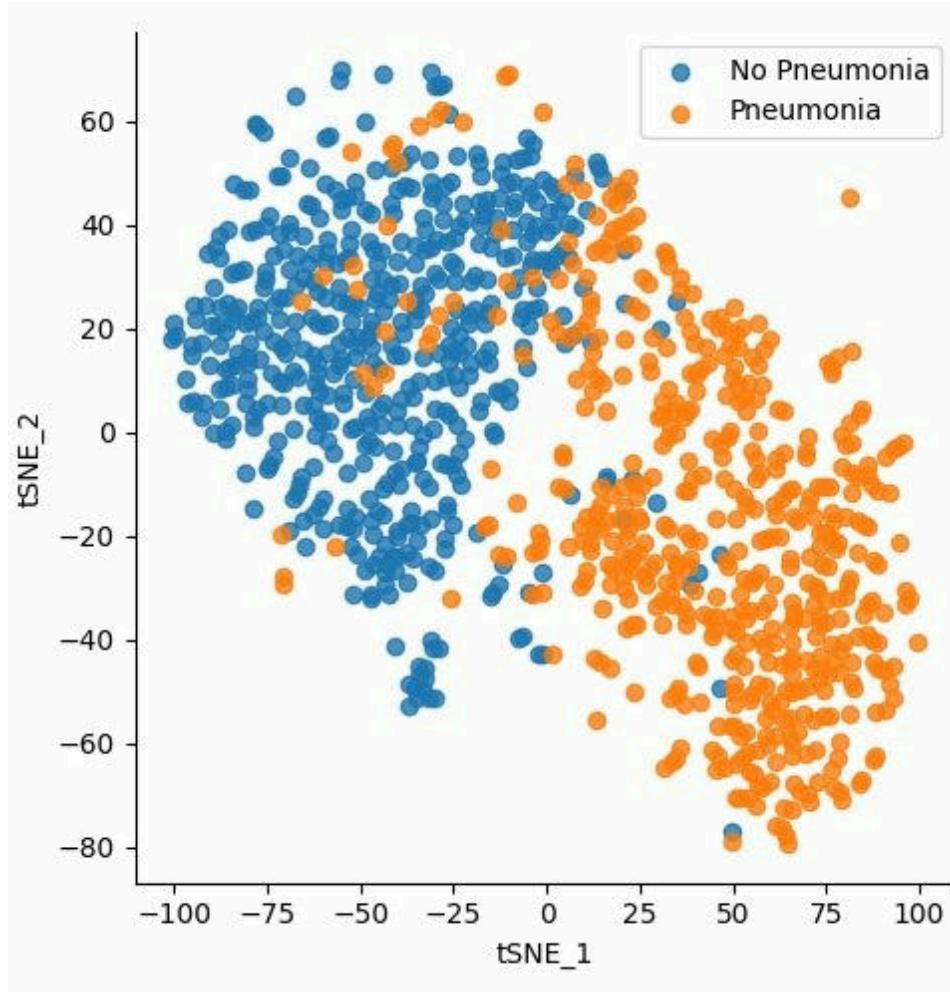
What is happening here? How can these models do this?

Autoencoders learn the structure of the dataset

Applications

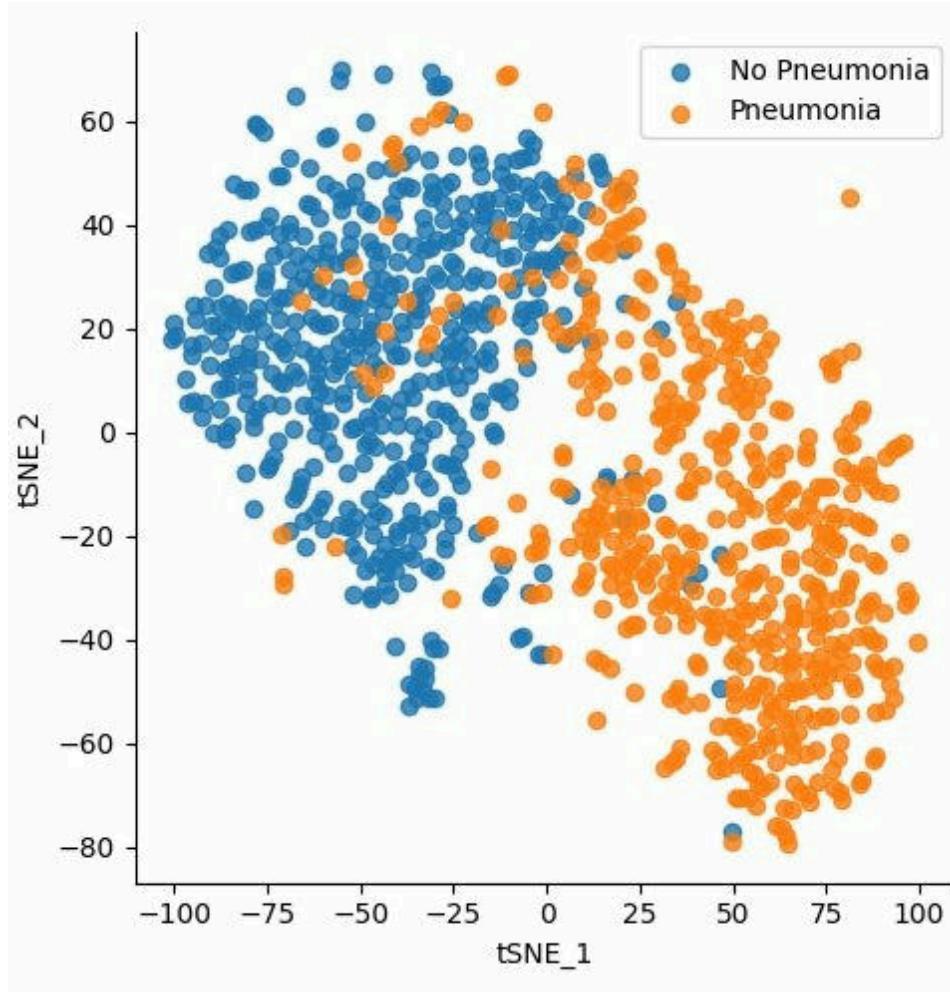


Applications



X : Pictures of human lungs

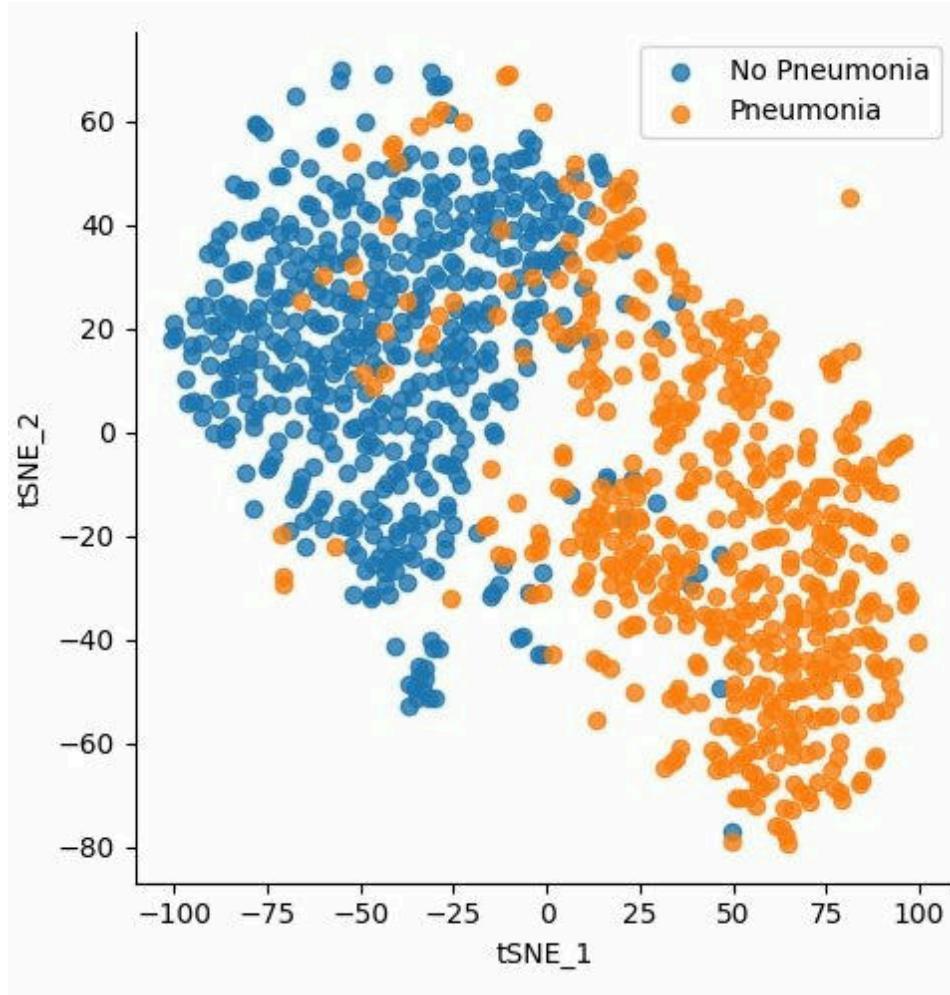
Applications



X : Pictures of human lungs

$$Z \in \mathbb{R}^2$$

Applications

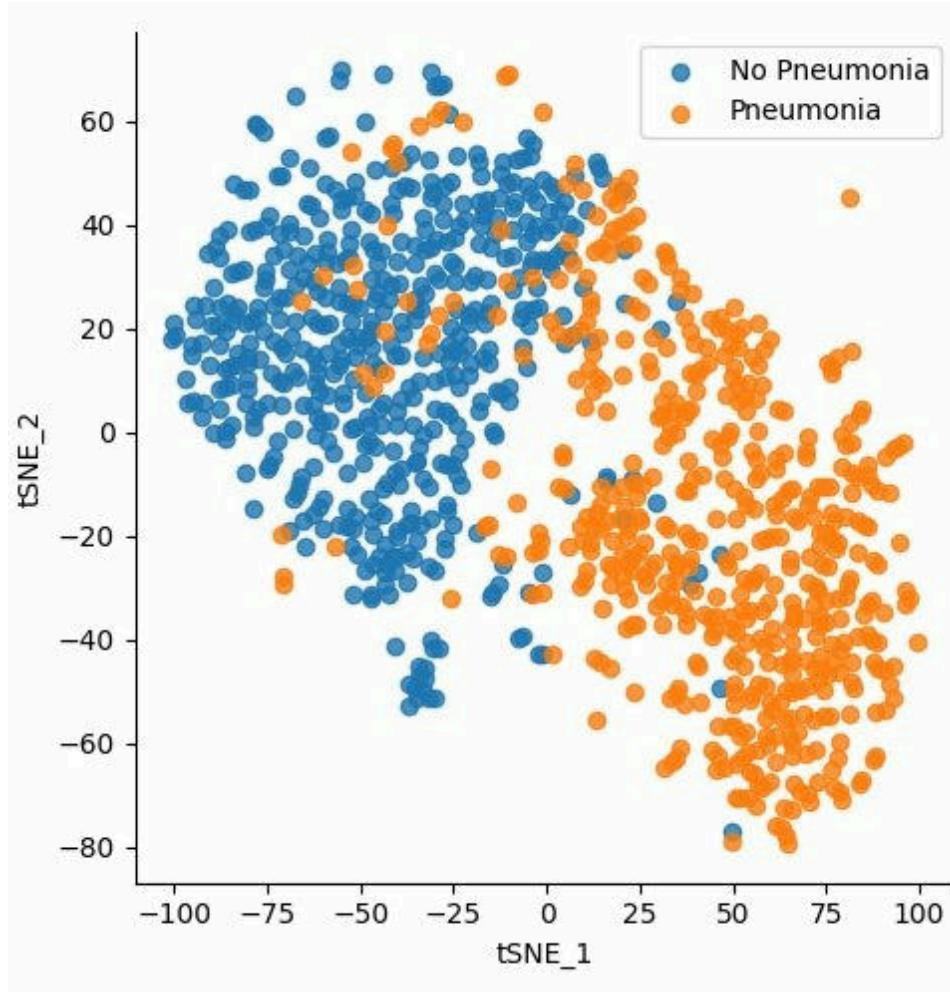


X : Pictures of human lungs

$$Z \in \mathbb{R}^2$$

Learns the structure of lungs from images

Applications



X : Pictures of human lungs

$$Z \in \mathbb{R}^2$$

Learns the structure of lungs from images

Differentiates sick and healthy lungs without being told

Applications

If the dataset is lung images, the model learns the structure of lungs

Applications

If the dataset is lung images, the model learns the structure of lungs

If the dataset is pictures from our world, then the autoencoders learn the structure of the world

Applications

If the dataset is lung images, the model learns the structure of lungs

If the dataset is pictures from our world, then the autoencoders learn the structure of the world

Nobody tells them what a dog or cat is

Applications

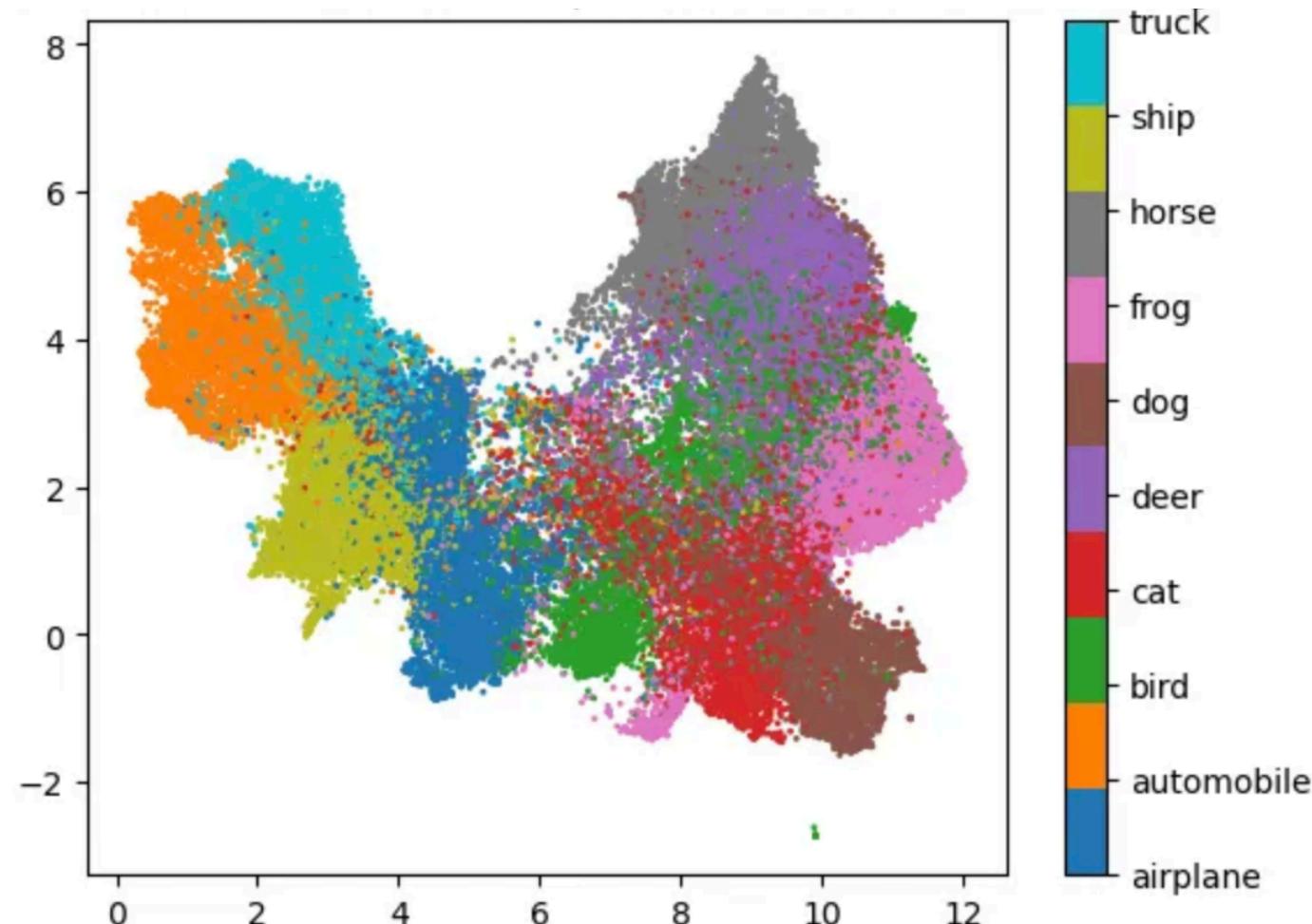
If the dataset is lung images, the model learns the structure of lungs

If the dataset is pictures from our world, then the autoencoders learn the structure of the world

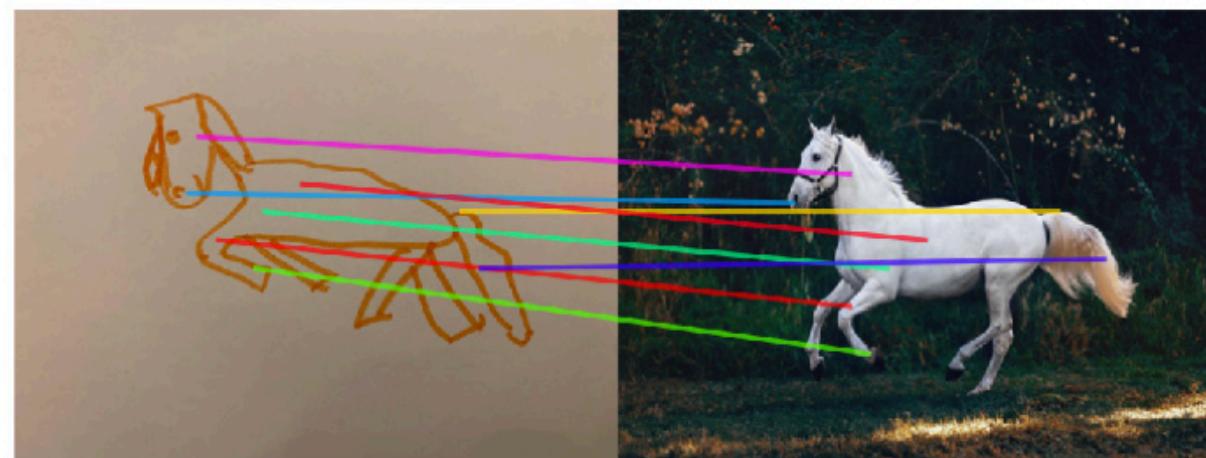
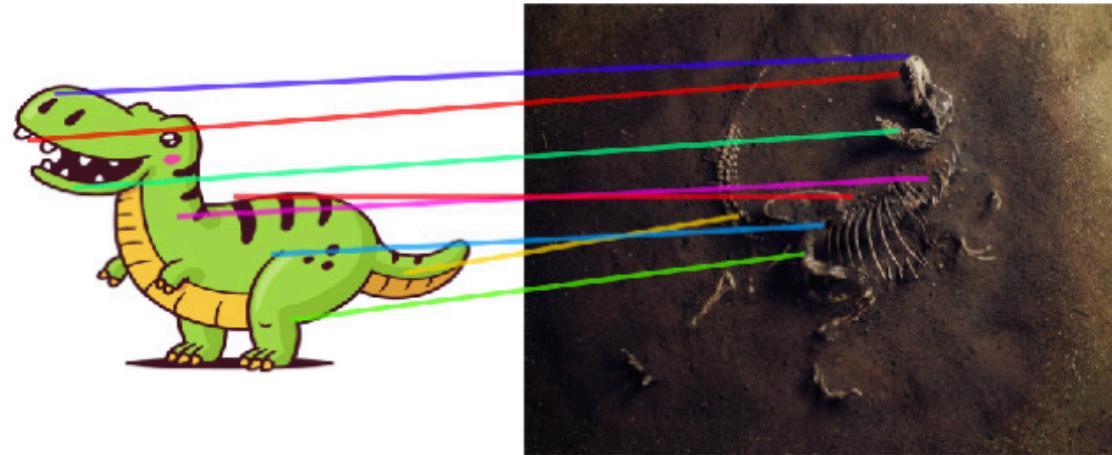
Nobody tells them what a dog or cat is

They learn that on their own

Applications



Applications



Applications

Some say “neural networks do not understand”, they just learn patterns

Applications

Some say “neural networks do not understand”, they just learn patterns

Humans are also pattern recognition machines

Applications

Some say “neural networks do not understand”, they just learn patterns

Humans are also pattern recognition machines



Applications

Some say “neural networks do not understand”, they just learn patterns

Humans are also pattern recognition machines



These networks can understand our world

Agenda

1. Review
2. Compression
3. Autoencoders
4. **Applications**
5. Variational Modeling
6. VAE Implementation
7. Coding

Agenda

1. Review
2. Compression
3. Autoencoders
4. Applications
5. **Variational Modeling**
6. VAE Implementation
7. Coding

Variational Modeling



Variational Modeling



These pictures were created by a **variational** autoencoder

Variational Modeling



These pictures were created by a **variational** autoencoder
But these people do not exist!

Variational Modeling

Autoencoders are useful for compression and denoising

Variational Modeling

Autoencoders are useful for compression and denoising

But we can also use them as **generative models**

Variational Modeling

Autoencoders are useful for compression and denoising

But we can also use them as **generative models**

A generative model learns the structure of data

Variational Modeling

Autoencoders are useful for compression and denoising

But we can also use them as **generative models**

A generative model learns the structure of data

Using this structure, it generates **new** data

Variational Modeling

Autoencoders are useful for compression and denoising

But we can also use them as **generative models**

A generative model learns the structure of data

Using this structure, it generates **new** data

- Train on face dataset, generate **new** pictures

Variational Modeling

Autoencoders are useful for compression and denoising

But we can also use them as **generative models**

A generative model learns the structure of data

Using this structure, it generates **new** data

- Train on face dataset, generate **new** pictures
- Train on book dataset, write a **new** book

Variational Modeling

Autoencoders are useful for compression and denoising

But we can also use them as **generative models**

A generative model learns the structure of data

Using this structure, it generates **new** data

- Train on face dataset, generate **new** pictures
- Train on book dataset, write a **new** book
- Train on protein dataset, create **new** proteins

Variational Modeling

Autoencoders are useful for compression and denoising

But we can also use them as **generative models**

A generative model learns the structure of data

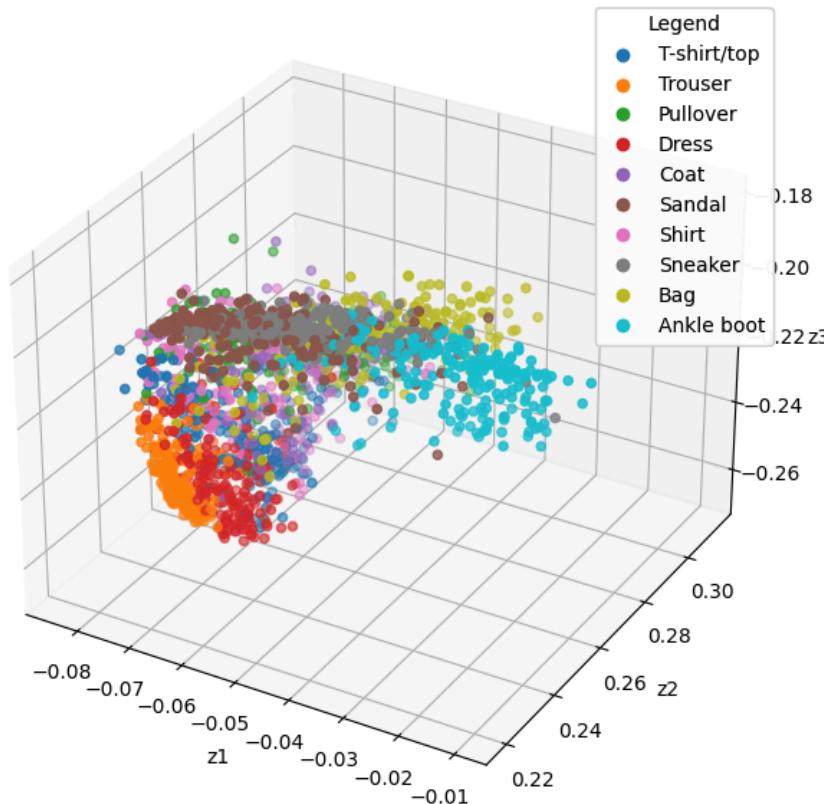
Using this structure, it generates **new** data

- Train on face dataset, generate **new** pictures
- Train on book dataset, write a **new** book
- Train on protein dataset, create **new** proteins

How does this work?

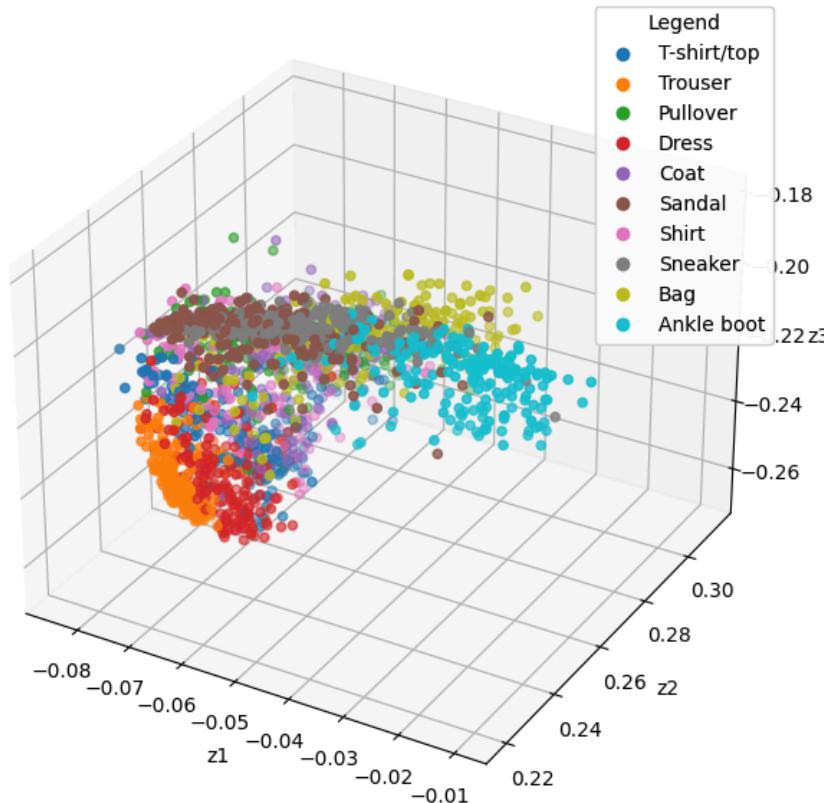
Variational Modeling

Latent space Z after training on the clothes dataset with $d_z = 3$

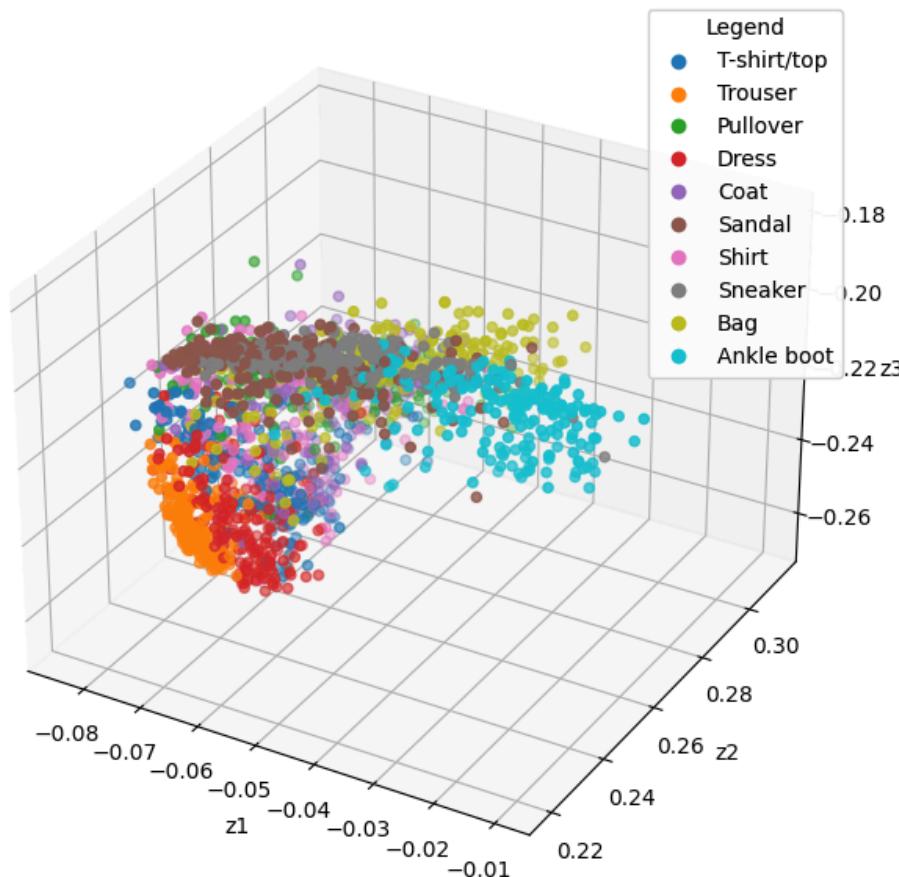


Variational Modeling

What happens if we decode a new point?

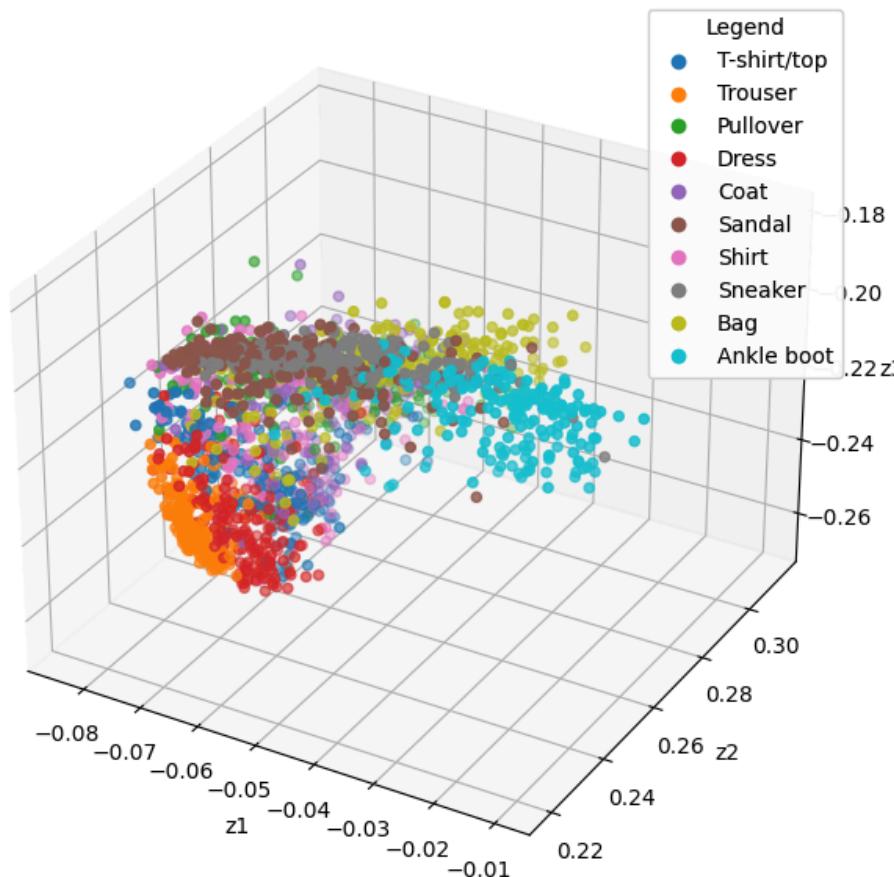


Variational Modeling



Autoencoder generative model:

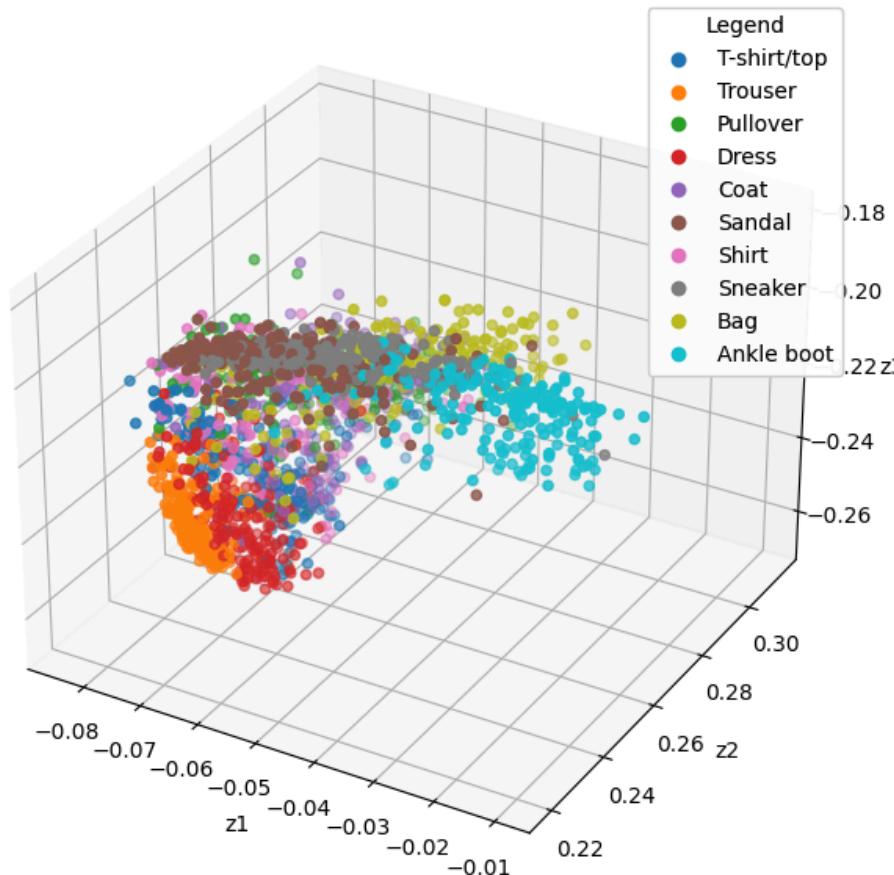
Variational Modeling



Autoencoder generative model:

Encode $\begin{bmatrix} \mathbf{x}_{[1]} \\ \vdots \\ \mathbf{x}_{[n]} \end{bmatrix}$ into $\begin{bmatrix} \mathbf{z}_{[1]} \\ \vdots \\ \mathbf{z}_{[n]} \end{bmatrix}$

Variational Modeling

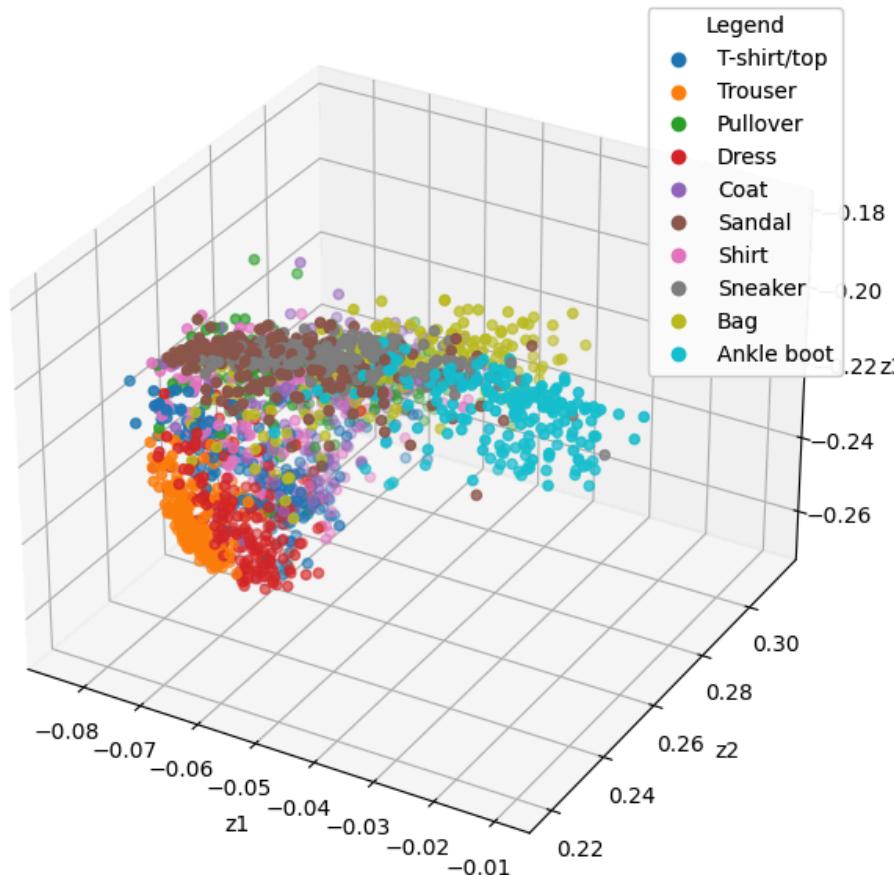


Autoencoder generative model:

Encode $\begin{bmatrix} \mathbf{x}_{[1]} \\ \vdots \\ \mathbf{x}_{[n]} \end{bmatrix}$ into $\begin{bmatrix} \mathbf{z}_{[1]} \\ \vdots \\ \mathbf{z}_{[n]} \end{bmatrix}$

Pick a point $\mathbf{z}_{[k]}$

Variational Modeling



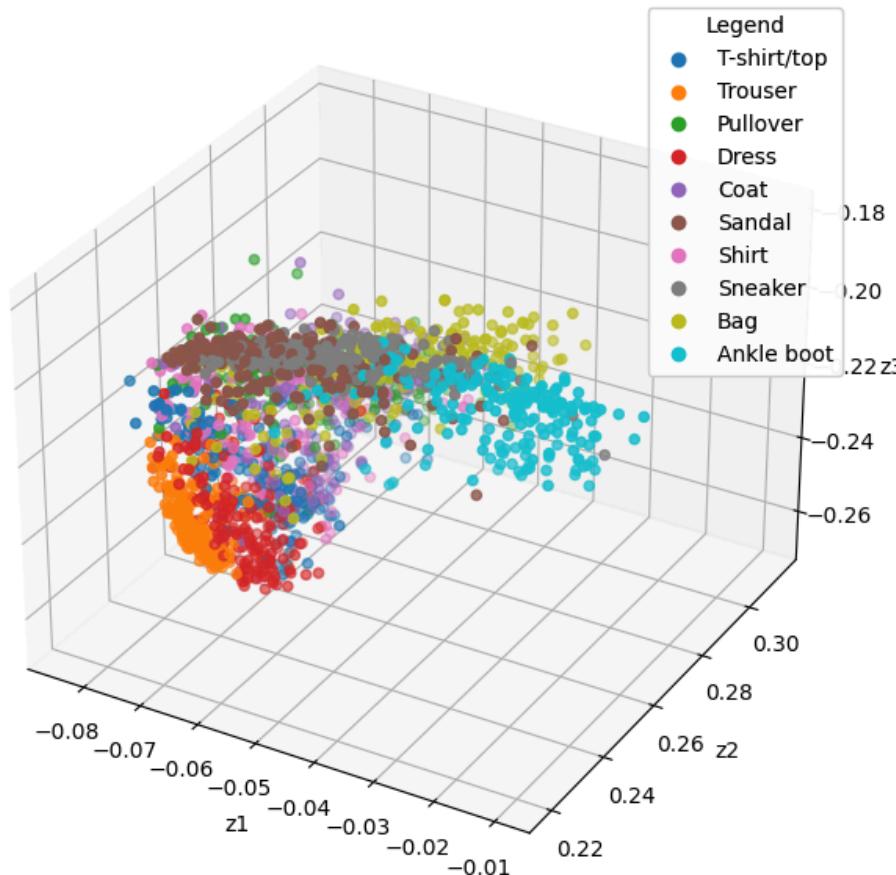
Autoencoder generative model:

Encode $\begin{bmatrix} \mathbf{x}_{[1]} \\ \vdots \\ \mathbf{x}_{[n]} \end{bmatrix}$ into $\begin{bmatrix} \mathbf{z}_{[1]} \\ \vdots \\ \mathbf{z}_{[n]} \end{bmatrix}$

Pick a point $\mathbf{z}_{[k]}$

Add some noise $\mathbf{z}_{\text{new}} = \mathbf{z}_{[k]} + \boldsymbol{\varepsilon}$

Variational Modeling



Autoencoder generative model:

$$\text{Encode } \begin{bmatrix} \mathbf{x}_{[1]} \\ \vdots \\ \mathbf{x}_{[n]} \end{bmatrix} \text{ into } \begin{bmatrix} \mathbf{z}_{[1]} \\ \vdots \\ \mathbf{z}_{[n]} \end{bmatrix}$$

Pick a point $\mathbf{z}_{[k]}$

Add some noise $\mathbf{z}_{\text{new}} = \mathbf{z}_{[k]} + \boldsymbol{\varepsilon}$

Decode \mathbf{z}_{new} into \mathbf{x}_{new}

Variational Modeling



Variational Modeling



$$f^{-1}(z_k + \epsilon, \theta_d)$$

But there is a problem, the **curse
of dimensionality**

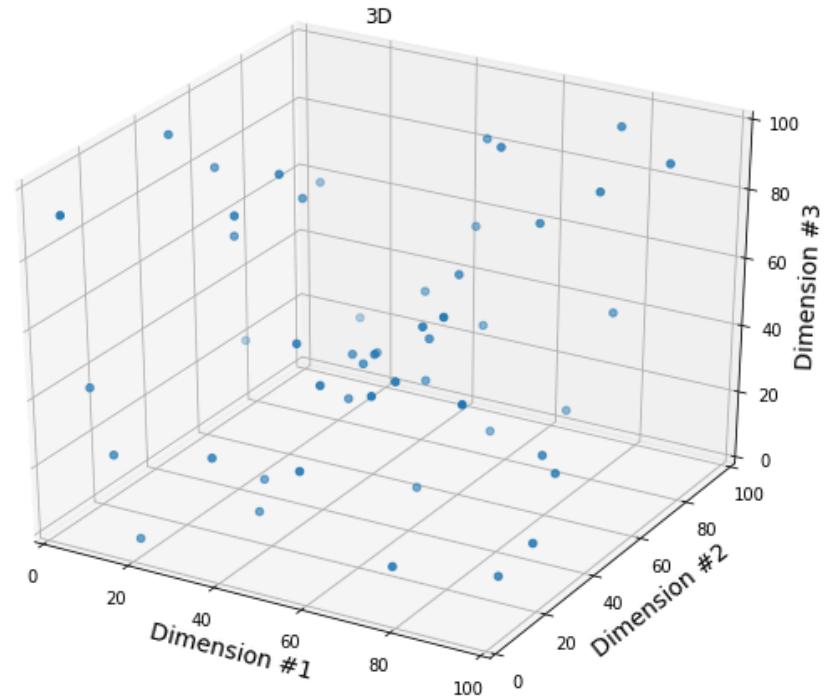
But there is a problem, the **curse
of dimensionality**

As d_z increases, points move
further and further apart

Variational Modeling

But there is a problem, the **curse of dimensionality**

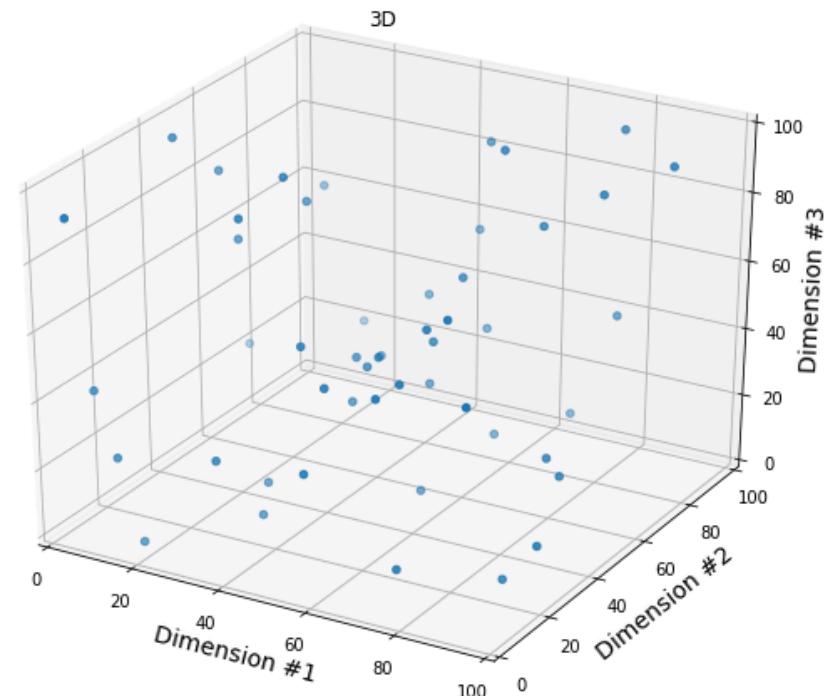
As d_z increases, points move further and further apart



Variational Modeling

But there is a problem, the **curse of dimensionality**

As d_z increases, points move further and further apart



$f^{-1}(z + \varepsilon)$ will produce either garbage, or z

Variational Modeling

Question: What can we do?

Variational Modeling

Question: What can we do?

Answer: Force the points to be close together!

Variational Modeling

Question: What can we do?

Answer: Force the points to be close together!

We will use a **variational autoencoder** (VAE)

Variational Modeling

VAE discovered by Diederik Kingma (also adam optimizer)

Variational Modeling

VAE discovered by Diederik Kingma (also adam optimizer)



Variational Modeling

Variational autoencoders (VAEs) do three things:

Variational Modeling

Variational autoencoders (VAEs) do three things:

1. Make it easy to sample random z

Variational Modeling

Variational autoencoders (VAEs) do three things:

1. Make it easy to sample random z
2. Keep all $z_{[1]}, \dots, z_{[n]}$ close together in a small region

Variational Modeling

Variational autoencoders (VAEs) do three things:

1. Make it easy to sample random z
2. Keep all $z_{[1]}, \dots, z_{[n]}$ close together in a small region
3. Ensure that $z + \epsilon$ is always meaningful

Variational Modeling

Variational autoencoders (VAEs) do three things:

1. Make it easy to sample random z
2. Keep all $z_{[1]}, \dots, z_{[n]}$ close together in a small region
3. Ensure that $z + \epsilon$ is always meaningful

How?

Variational Modeling

Variational autoencoders (VAEs) do three things:

1. Make it easy to sample random z
2. Keep all $z_{[1]}, \dots, z_{[n]}$ close together in a small region
3. Ensure that $z + \epsilon$ is always meaningful

How?

Make $z_{[1]}, \dots, z_{[n]}$ normally distributed

Variational Modeling

Variational autoencoders (VAEs) do three things:

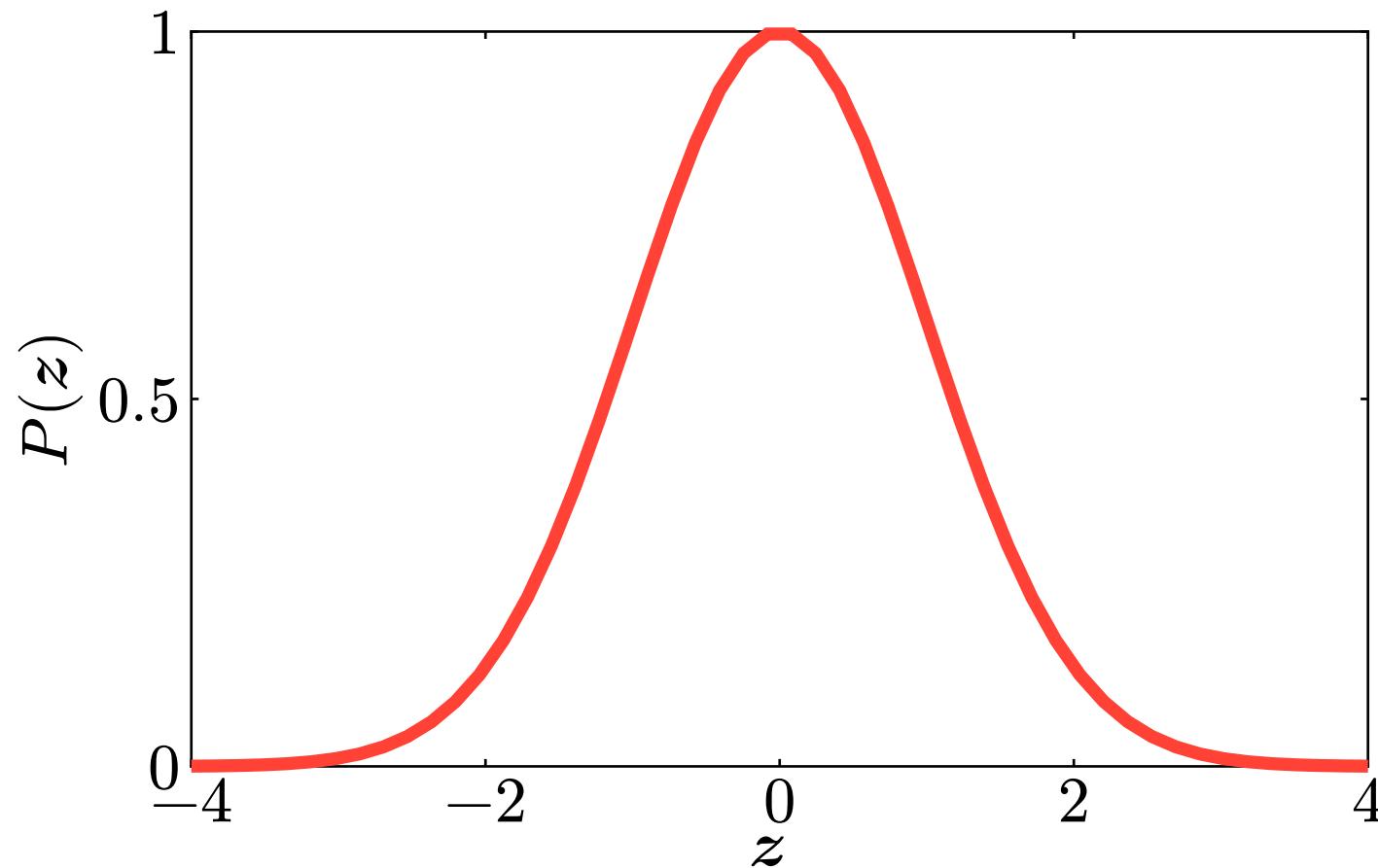
1. Make it easy to sample random z
2. Keep all $z_{[1]}, \dots, z_{[n]}$ close together in a small region
3. Ensure that $z + \epsilon$ is always meaningful

How?

Make $z_{[1]}, \dots, z_{[n]}$ normally distributed

$$z \sim \mathcal{N}(\mu, \sigma), \quad \mu = 0, \sigma = 1$$

Variational Modeling



Variational Modeling

If $z_{[1]}, \dots, z_{[n]}$ are distributed following $\mathcal{N}(0, 1)$:

Variational Modeling

If $z_{[1]}, \dots, z_{[n]}$ are distributed following $\mathcal{N}(0, 1)$:

1. 99.7% of $z_{[1]}, \dots, z_{[n]}$ lie within $3\sigma = [-3, 3]$

Variational Modeling

If $z_{[1]}, \dots, z_{[n]}$ are distributed following $\mathcal{N}(0, 1)$:

1. 99.7% of $z_{[1]}, \dots, z_{[n]}$ lie within $3\sigma = [-3, 3]$
2. Make it easy to generate new z , just sample $z \sim \mathcal{N}(0, 1)$

Variational Modeling

So how do we ensure that $z_{[1]}, \dots, z_{[n]}$ are normally distributed?

Variational Modeling

So how do we ensure that $z_{[1]}, \dots, z_{[n]}$ are normally distributed?

We have to remember conditional probabilities

Variational Modeling

So how do we ensure that $z_{[1]}, \dots, z_{[n]}$ are normally distributed?

We have to remember conditional probabilities

$P(\text{rain} \mid \text{cloud})$ = Probability of rain, given that it is cloudy

Variational Modeling

So how do we ensure that $z_{[1]}, \dots, z_{[n]}$ are normally distributed?

We have to remember conditional probabilities

$P(\text{rain} \mid \text{cloud})$ = Probability of rain, given that it is cloudy

Variational Modeling

Key idea 1: We want to model the distribution over the dataset X

$$P(x; \theta), \quad x \sim X$$

Variational Modeling

Key idea 1: We want to model the distribution over the dataset X

$$P(x; \theta), \quad x \sim X$$

We want to learn θ that best models the distribution of possible faces

Variational Modeling

Key idea 1: We want to model the distribution over the dataset X

$$P(x; \theta), \quad x \sim X$$

We want to learn θ that best models the distribution of possible faces

Large $P(x; \theta)$



$P(x; \theta) \approx 0$



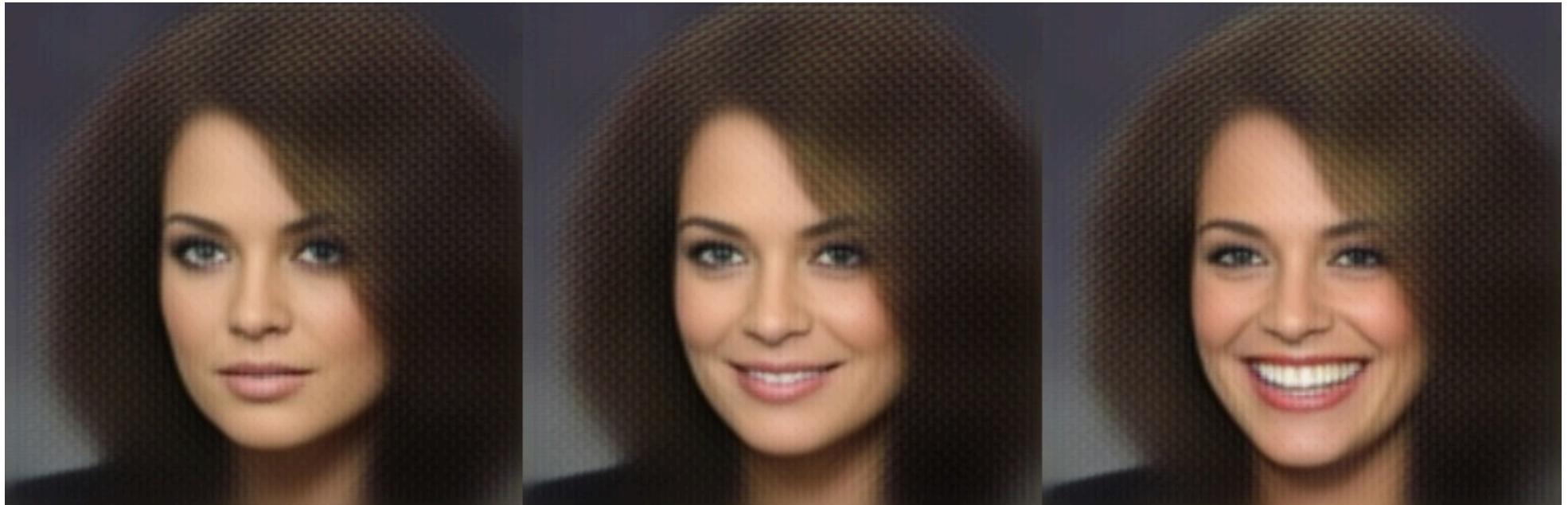
Variational Modeling

Key idea 2: There is some latent variable z which generates data x

Variational Modeling

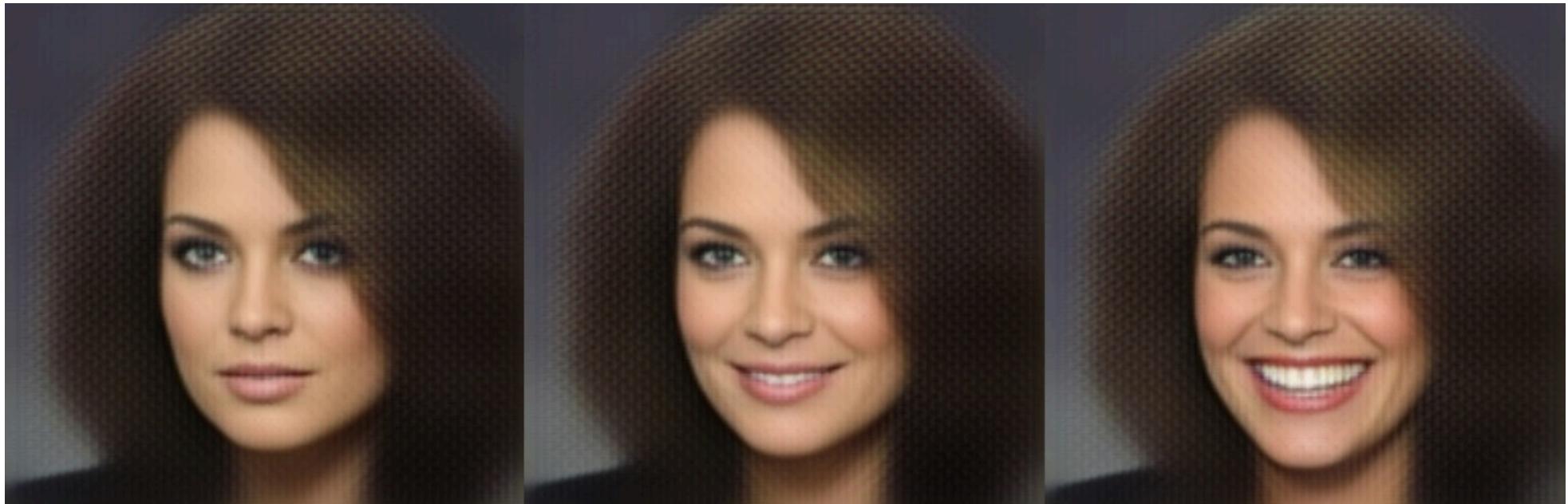
Key idea 2: There is some latent variable z which generates data x

$x :$



$z : [\text{woman brown hair (frown} \mid \text{smile})]$

Variational Modeling

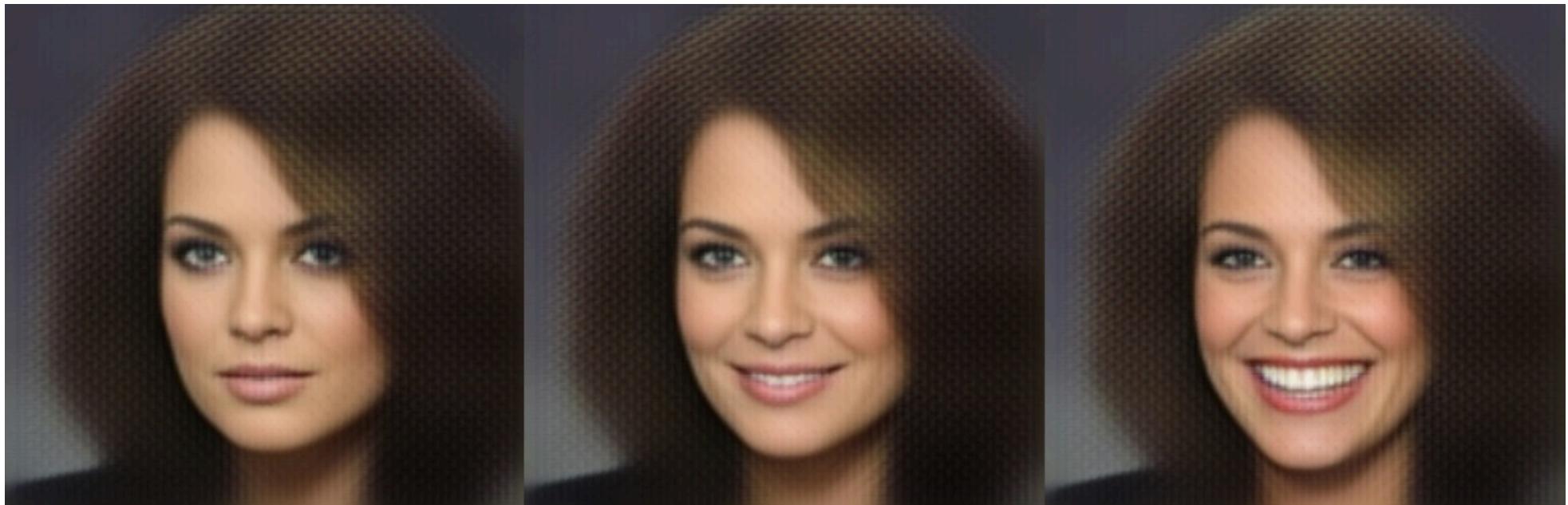


Variational Modeling



Network can only see x , it cannot directly observe z

Variational Modeling



Network can only see x , it cannot directly observe z

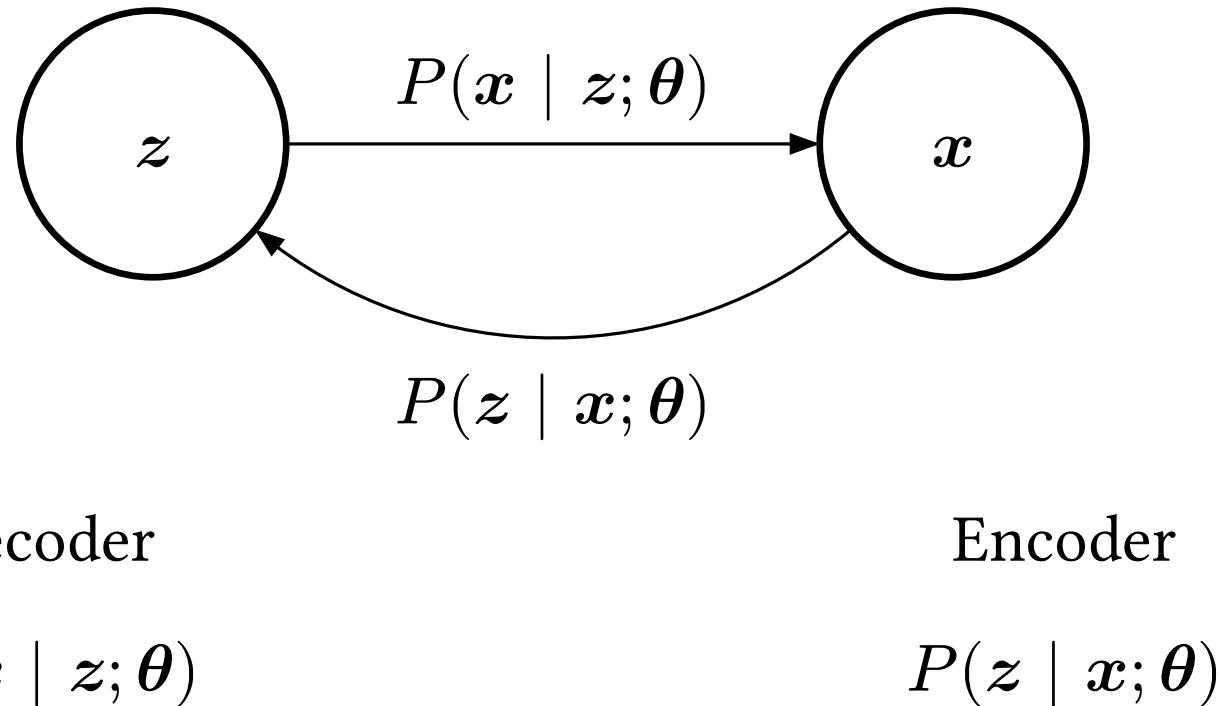
Given x , find the probability that the person is smiling $P(z \mid x; \theta)$

Variational Modeling

We cast the autoencoding task as a **variational inference** problem

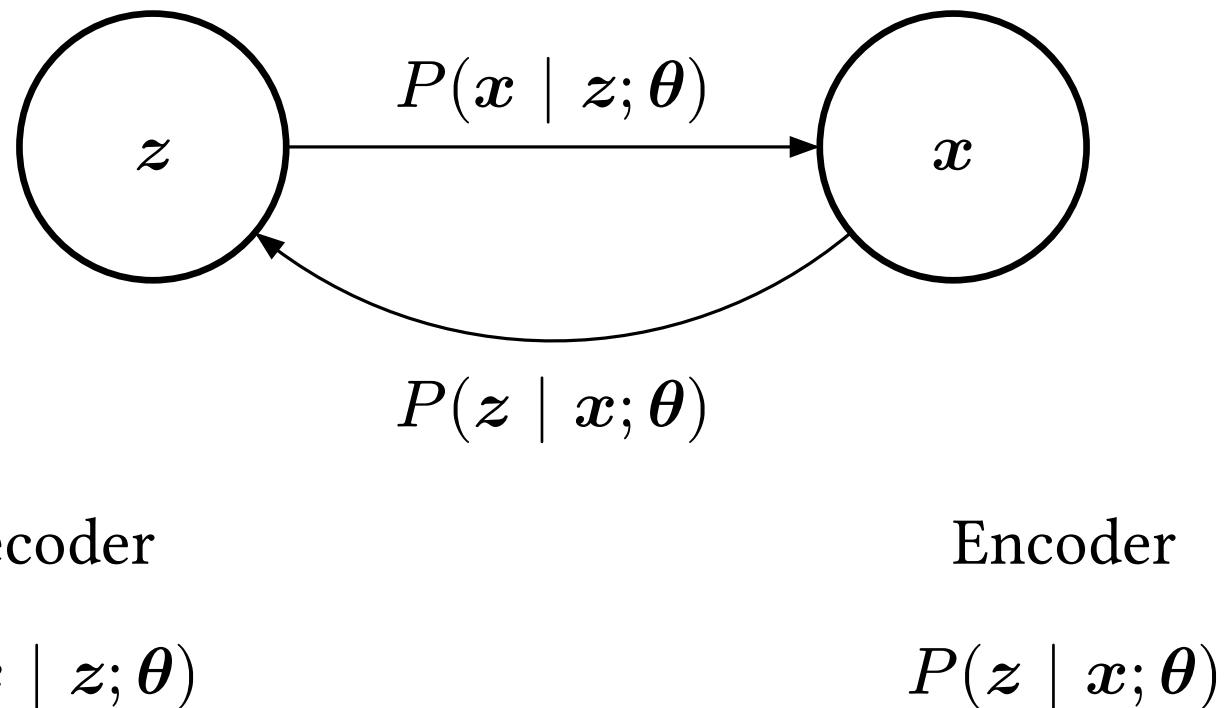
Variational Modeling

We cast the autoencoding task as a **variational inference** problem



Variational Modeling

We cast the autoencoding task as a **variational inference** problem



We want to learn both the encoder and decoder: $P(z, x; \theta)$

Variational Modeling

$$P(z, x; \theta) = P(x \mid z; \theta) P(z; \theta)$$

Variational Modeling

$$P(z, x; \theta) = P(x | z; \theta) P(z; \theta)$$

We can choose any distribution for $P(z)$

Variational Modeling

$$P(z, x; \theta) = P(x | z; \theta) P(z; \theta)$$

We can choose any distribution for $P(z)$

$$P(z) = \mathcal{N}(\mathbf{0}, \mathbf{1})$$

Variational Modeling

$$P(z, x; \theta) = P(x | z; \theta) P(z; \theta)$$

We can choose any distribution for $P(z)$

$$P(z) = \mathcal{N}(\mathbf{0}, \mathbf{1})$$

We can generate all possible x by sampling $z \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$

Variational Modeling

$$P(z, x; \theta) = P(x | z; \theta) P(z; \theta)$$

We can choose any distribution for $P(z)$

$$P(z) = \mathcal{N}(\mathbf{0}, \mathbf{1})$$

We can generate all possible x by sampling $z \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$

We can randomly generate z , which we can decode into new x !

Variational Modeling

Now, all we must do is find θ that best explains the dataset distribution

Variational Modeling

Now, all we must do is find θ that best explains the dataset distribution

Learned distribution $P(x; \theta)$ to be close to dataset $P(x)$, $x \sim X$

Variational Modeling

Now, all we must do is find θ that best explains the dataset distribution

Learned distribution $P(x; \theta)$ to be close to dataset $P(x)$, $x \sim X$

We need some error function between $P(x; \theta)$ and $P(x)$

Variational Modeling

Now, all we must do is find θ that best explains the dataset distribution

Learned distribution $P(x; \theta)$ to be close to dataset $P(x)$, $x \sim X$

We need some error function between $P(x; \theta)$ and $P(x)$

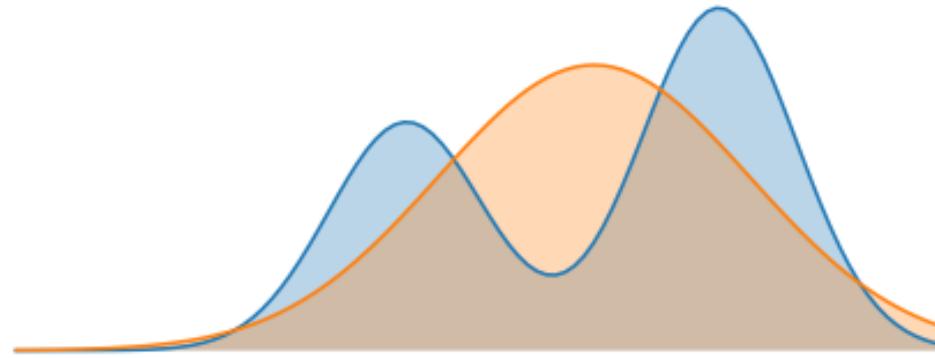
Question: How do we measure the distance between probability distributions?

Variational Modeling

Answer: KL divergence

Variational Modeling

Answer: KL divergence



$$\text{KL}(P, Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

Variational Modeling

Learn the parameters for our model

Variational Modeling

Learn the parameters for our model

$$\arg \min_{\theta} \text{KL}(P(x), P(x; \theta))$$

Variational Modeling

Learn the parameters for our model

$$\arg \min_{\theta} \text{KL}(P(x), P(x; \theta))$$

Unfortunately, this objective is intractable to optimize

Variational Modeling

Learn the parameters for our model

$$\arg \min_{\theta} \text{KL}(P(x), P(x; \theta))$$

Unfortunately, this objective is intractable to optimize

The paper provides surrogate objective

$$\arg \min_{\theta} \left[-\log P(x \mid z; \theta) + \frac{1}{2} \text{KL}(P(z \mid x; \theta), P(z)) \right]$$

Variational Modeling

Learn the parameters for our model

$$\arg \min_{\theta} \text{KL}(P(x), P(x; \theta))$$

Unfortunately, this objective is intractable to optimize

The paper provides surrogate objective

$$\arg \min_{\theta} \left[-\log P(x \mid z; \theta) + \frac{1}{2} \text{KL}(P(z \mid x; \theta), P(z)) \right]$$

We call this the **Evidence Lower Bound Objective** (ELBO)

Variational Modeling

$$\arg \min_{\theta} \left[-\log P(x \mid z; \theta) + \frac{1}{2} \text{KL}(P(z \mid x; \theta), P(z)) \right]$$

Variational Modeling

$$\arg \min_{\theta} \left[-\log P(x \mid z; \theta) + \frac{1}{2} \text{KL}(P(z \mid x; \theta), P(z)) \right]$$

How is this ELBO helpful?

Variational Modeling

$$\arg \min_{\theta} \left[-\log P(x \mid z; \theta) + \frac{1}{2} \text{KL}(P(z \mid x; \theta), P(z)) \right]$$

How is this ELBO helpful?

Decoder

$$P(x \mid z; \theta)$$

Encoder

$$P(z \mid x; \theta)$$

Prior

$$P(z) = \mathcal{N}(\mathbf{0}, \mathbf{1})$$

Variational Modeling

$$\arg \min_{\theta} \left[-\log P(x \mid z; \theta) + \frac{1}{2} \text{KL}(P(z \mid x; \theta), P(z)) \right]$$

How is this ELBO helpful?

Decoder

$$P(x \mid z; \theta)$$

Encoder

$$P(z \mid x; \theta)$$

Prior

$$P(z) = \mathcal{N}(\mathbf{0}, \mathbf{1})$$

$$\arg \min_{\theta} \left[\underbrace{-\log P(x \mid z; \theta)}_{\text{Reconstruction error}} + \frac{1}{2} \underbrace{\text{KL}(P(z \mid x; \theta), P(z))}_{\text{Constrain latent}} \right]$$

Variational Modeling

$$\arg \min_{\theta} \left[-\log P(x \mid z; \theta) + \frac{1}{2} \text{KL}(P(z \mid x; \theta), P(z)) \right]$$

How is this ELBO helpful?

Decoder	Encoder	Prior
$P(x \mid z; \theta)$	$P(z \mid x; \theta)$	$P(z) = \mathcal{N}(\mathbf{0}, \mathbf{1})$

$$\arg \min_{\theta} \left[\underbrace{-\log P(x \mid z; \theta)}_{\text{Reconstruction error}} + \frac{1}{2} \underbrace{\text{KL}(P(z \mid x; \theta), P(z))}_{\text{Constrain latent}} \right]$$

Now we know how to train our autoencoder!

Agenda

1. Review
2. Compression
3. Autoencoders
4. Applications
5. **Variational Modeling**
6. VAE Implementation
7. Coding

Agenda

1. Review
2. Compression
3. Autoencoders
4. Applications
5. Variational Modeling
6. **VAE Implementation**
7. Coding

VAE Implementation

How do we implement f (i.e., $P(z \mid x; \theta)$)?

VAE Implementation

How do we implement f (i.e., $P(z | x; \theta)$)?

$$f : X \times \Theta \mapsto \Delta Z$$

VAE Implementation

How do we implement f (i.e., $P(z | x; \theta)$)?

$$f : X \times \Theta \mapsto \Delta Z$$

Normal distribution has a mean $\mu \in \mathbb{R}$ and standard deviation $\sigma \in \mathbb{R}_+$

VAE Implementation

How do we implement f (i.e., $P(z | x; \theta)$)?

$$f : X \times \Theta \mapsto \Delta Z$$

Normal distribution has a mean $\mu \in \mathbb{R}$ and standard deviation $\sigma \in \mathbb{R}_+$

Our encoder should output d_z means and d_z standard deviations

VAE Implementation

How do we implement f (i.e., $P(z | x; \theta)$)?

$$f : X \times \Theta \mapsto \Delta Z$$

Normal distribution has a mean $\mu \in \mathbb{R}$ and standard deviation $\sigma \in \mathbb{R}_+$

Our encoder should output d_z means and d_z standard deviations

$$f : X \times \Theta \mapsto \mathbb{R}^{d_z} \times \mathbb{R}_+^{d_z}$$

VAE Implementation

```
core = nn.Sequential(...)  
mu_layer = nn.Linear(d_h, d_z)  
# Neural networks output real numbers  
# But sigma must be positive  
# Output log sigma, because e^(sigma) is always positive  
log_sigma_layer = nn.Linear(d_h, d_z)  
# Alternatively, one sigma for all data  
log_sigma = jnp.ones((d_z,))  
  
tmp = core(x)  
mu = mu_layer(tmp)  
log_sigma = log_sigma_layer(tmp)  
distribution = (mu, exp(sigma))
```

VAE Implementation

We covered the encoder

$$f : X \times \Theta \mapsto \Delta Z$$

We can use the same decoder as a standard autoencoder

VAE Implementation

We covered the encoder

$$f : X \times \Theta \mapsto \Delta Z$$

We can use the same decoder as a standard autoencoder

$$f^{-1} : Z \times \Theta \mapsto X$$

Question: Any issues?

VAE Implementation

We covered the encoder

$$f : X \times \Theta \mapsto \Delta Z$$

We can use the same decoder as a standard autoencoder

$$f^{-1} : Z \times \Theta \mapsto X$$

Question: Any issues?

Answer: Encoder outputs a distribution ΔZ but decoder input is Z

VAE Implementation

We can sample from the distribution

$$\mu, \sigma = f(x, \theta_e)$$

$$z \sim \mathcal{N}(\mu, \sigma)$$

VAE Implementation

We can sample from the distribution

$$\mu, \sigma = f(x, \theta_e)$$

$$z \sim \mathcal{N}(\mu, \sigma)$$

But there is a problem! Sampling is not differentiable

VAE Implementation

We can sample from the distribution

$$\mu, \sigma = f(x, \theta_e)$$

$$z \sim \mathcal{N}(\mu, \sigma)$$

But there is a problem! Sampling is not differentiable

Question: Why does this matter?

VAE Implementation

We can sample from the distribution

$$\mu, \sigma = f(x, \theta_e)$$

$$z \sim \mathcal{N}(\mu, \sigma)$$

But there is a problem! Sampling is not differentiable

Question: Why does this matter?

Answer: Must be differentiable for gradient descent

VAE Implementation

VAE paper proposes the **reparameterization trick**

VAE Implementation

VAE paper proposes the **reparameterization trick**

$$z \sim \mathcal{N}(\mu, \sigma)$$

$$z = \mu + \sigma \odot \varepsilon \quad \varepsilon \sim \mathcal{N}(0, 1)$$

VAE Implementation

VAE paper proposes the **reparameterization trick**

$$z \sim \mathcal{N}(\mu, \sigma)$$

$$z = \mu + \sigma \odot \varepsilon \quad \varepsilon \sim \mathcal{N}(0, 1)$$

Gradient can flow through μ, σ

VAE Implementation

VAE paper proposes the **reparameterization trick**

$$z \sim \mathcal{N}(\mu, \sigma)$$

$$z = \mu + \sigma \odot \varepsilon \quad \varepsilon \sim \mathcal{N}(0, 1)$$

Gradient can flow through μ, σ

We can sample and use gradient descent

VAE Implementation

VAE paper proposes the **reparameterization trick**

$$z \sim \mathcal{N}(\mu, \sigma)$$

$$z = \mu + \sigma \odot \varepsilon \quad \varepsilon \sim \mathcal{N}(0, 1)$$

Gradient can flow through μ, σ

We can sample and use gradient descent

This trick only works with certain distributions

VAE Implementation

Put it all together

VAE Implementation

Put it all together

Step 1: Encode the input to a normal distribution

$$\mu, \sigma = f(x, \theta_e)$$

VAE Implementation

Put it all together

Step 1: Encode the input to a normal distribution

$$\mu, \sigma = f(x, \theta_e)$$

Step 2: Generate a sample from distribution

$$z = \mu + \sigma \odot \varepsilon$$

VAE Implementation

Put it all together

Step 1: Encode the input to a normal distribution

$$\mu, \sigma = f(x, \theta_e)$$

Step 2: Generate a sample from distribution

$$z = \mu + \sigma \odot \varepsilon$$

Step 3: Decode the sample

$$x = f^{-1}(z, \theta_d)$$

VAE Implementation

One last thing, implement the loss function

VAE Implementation

One last thing, implement the loss function

Decoder	Encoder	Prior
$P(x \mid z; \theta)$	$P(z \mid x; \theta)$	$P(z) = \mathcal{N}(0, 1)$

VAE Implementation

One last thing, implement the loss function

Decoder

$$P(x \mid z; \theta)$$

Encoder

$$P(z \mid x; \theta)$$

Prior

$$P(z) = \mathcal{N}(0, 1)$$

$$\mathcal{L}(x, \theta) = \arg \min_{\theta} \left[\underbrace{-\log P(x \mid z; \theta)}_{\text{Reconstruction error}} + \underbrace{\frac{1}{2} \text{KL}(P(z \mid x; \theta), P(z))}_{\text{Constrain latent}} \right]$$

VAE Implementation

One last thing, implement the loss function

Decoder	Encoder	Prior
$P(x \mid z; \theta)$	$P(z \mid x; \theta)$	$P(z) = \mathcal{N}(0, 1)$

$$\mathcal{L}(x, \theta) = \arg \min_{\theta} \left[\underbrace{-\log P(x \mid z; \theta)}_{\text{Reconstruction error}} + \underbrace{\frac{1}{2} \text{KL}(P(z \mid x; \theta), P(z))}_{\text{Constrain latent}} \right]$$

Start with the KL term first

VAE Implementation

$$\mathcal{L}(x, \theta) = \arg \min_{\theta} \left[-\log P(x \mid z; \theta) + \frac{1}{2} \text{KL}(P(z \mid x; \theta), P(z)) \right]$$

VAE Implementation

$$\mathcal{L}(x, \theta) = \arg \min_{\theta} \left[-\log P(x \mid z; \theta) + \frac{1}{2} \text{KL}(P(z \mid x; \theta), P(z)) \right]$$

First, rewrite KL term using our encoder f

VAE Implementation

$$\mathcal{L}(x, \theta) = \arg \min_{\theta} \left[-\log P(x | z; \theta) + \frac{1}{2} \text{KL}(P(z | x; \theta), P(z)) \right]$$

First, rewrite KL term using our encoder f

$$\mathcal{L}(x, \theta) = \arg \min_{\theta} \left[-\log P(x | z) + \frac{1}{2} \text{KL}(f(x, \theta_e), P(z)) \right]$$

VAE Implementation

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \left[-\log P(\mathbf{x} \mid \mathbf{z}; \boldsymbol{\theta}) + \frac{1}{2} \text{KL}(P(\mathbf{z} \mid \mathbf{x}; \boldsymbol{\theta}), P(\mathbf{z})) \right]$$

First, rewrite KL term using our encoder f

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \left[-\log P(\mathbf{x} \mid \mathbf{z}) + \frac{1}{2} \text{KL}(f(\mathbf{x}, \boldsymbol{\theta}_e), P(\mathbf{z})) \right]$$

$P(\mathbf{z})$ and $f(\mathbf{x}, \boldsymbol{\theta}_e)$ are Gaussian, we can simplify KL term

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\theta}) = \underbrace{\log P(\mathbf{x} \mid \mathbf{z})}_{\text{Reconstruction error}} - \left(\sum_{j=1}^{d_z} \mu_j^2 + \sigma_j^2 - \log(\sigma^2) - 1 \right)$$

VAE Implementation

$$\mathcal{L}(\mathbf{x}, \theta) = \underbrace{\log P(\mathbf{x} \mid \mathbf{z})}_{\text{Reconstruction error}} - \left(\sum_{j=1}^{d_z} \mu_j^2 + \sigma_j^2 - \log(\sigma^2) - 1 \right)$$

VAE Implementation

$$\mathcal{L}(\mathbf{x}, \theta) = \underbrace{\log P(\mathbf{x} \mid \mathbf{z})}_{\text{Reconstruction error}} - \left(\sum_{j=1}^{d_z} \mu_j^2 + \sigma_j^2 - \log(\sigma^2) - 1 \right)$$

Next, plug in square error for reconstruction error

VAE Implementation

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\theta}) = \underbrace{\log P(\mathbf{x} \mid \mathbf{z})}_{\text{Reconstruction error}} - \left(\sum_{j=1}^{d_z} \mu_j^2 + \sigma_j^2 - \log(\sigma^2) - 1 \right)$$

Next, plug in square error for reconstruction error

$$= \sum_{j=1}^{d_z} \left(x_j - f^{-1}(f(\mathbf{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)_j \right)^2 - \left(\sum_{j=1}^{d_z} \mu_j^2 + \sigma_j^2 - \log(\sigma_j^2) - 1 \right)$$

VAE Implementation

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\theta}) = \underbrace{\log P(\mathbf{x} \mid \mathbf{z})}_{\text{Reconstruction error}} - \left(\sum_{j=1}^{d_z} \mu_j^2 + \sigma_j^2 - \log(\sigma^2) - 1 \right)$$

Next, plug in square error for reconstruction error

$$= \sum_{j=1}^{d_z} \left(x_j - f^{-1}(f(\mathbf{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)_j \right)^2 - \left(\sum_{j=1}^{d_z} \mu_j^2 + \sigma_j^2 - \log(\sigma_j^2) - 1 \right)$$

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^{d_z} \left(x_j - f^{-1}(f(\mathbf{x}, \boldsymbol{\theta}_e), \boldsymbol{\theta}_d)_j \right)^2 - \left(\sum_{j=1}^{d_z} \mu_j^2 + \sigma_j^2 - \log(\sigma_j^2) - 1 \right)$$

VAE Implementation

$$\mathcal{L}(x, \theta) = \sum_{j=1}^{d_z} \left(x_j - f^{-1}(f(x, \theta_e), \theta_d)_j \right)^2 - \left(\sum_{j=1}^{d_z} \mu_j^2 + \sigma_j^2 - \log(\sigma_j^2) - 1 \right)$$

VAE Implementation

$$\mathcal{L}(x, \theta) = \sum_{j=1}^{d_z} \left(x_j - f^{-1}(f(x, \theta_e), \theta_d)_j \right)^2 - \left(\sum_{j=1}^{d_z} \mu_j^2 + \sigma_j^2 - \log(\sigma_j^2) - 1 \right)$$

Finally, define over the entire dataset

$$\mathcal{L}(X, \theta) = \sum_{i=1}^n \sum_{j=1}^{d_z} \left(x_{[i],j} - f^{-1}(f(x_{[i]}, \theta_e), \theta_d)_j \right)^2 - \left(\sum_{i=1}^n \sum_{j=1}^{d_z} \mu_{[i],j}^2 + \sigma_{[i],j}^2 - \log(\sigma_{[i],j}^2) - 1 \right)$$

VAE Implementation

$$\begin{aligned}\mathcal{L}(\mathbf{X}, \boldsymbol{\theta}) = & \sum_{i=1}^n \sum_{j=1}^{d_z} \left(x_{[i],j} - f^{-1} \left(f \left(\mathbf{x}_{[i]}, \boldsymbol{\theta}_e \right), \boldsymbol{\theta}_d \right)_j \right)^2 - \\ & \left(\sum_{i=1}^n \sum_{j=1}^{d_z} \mu_{[i],j}^2 + \sigma_{[i],j}^2 - \log(\sigma_{[i],j}^2) - 1 \right)\end{aligned}$$

Scale of two terms can vary, we do not want one term to dominate

VAE Implementation

Paper suggests using minibatch size m and dataset size n

VAE Implementation

Paper suggests using minibatch size m and dataset size n

$$\mathcal{L}(\mathbf{X}, \boldsymbol{\theta}) = \frac{m}{n} \sum_{i=1}^n \sum_{j=1}^{d_z} \left(x_{[i],j} - f^{-1}\left(f\left(\mathbf{x}_{[i]}, \boldsymbol{\theta}_e\right), \boldsymbol{\theta}_d\right)_j \right)^2 - \\ \left(\sum_{i=1}^n \sum_{j=1}^{d_z} \mu_{[i],j}^2 + \sigma_{[i],j}^2 - \log(\sigma_{[i],j}^2) - 1 \right)$$

VAE Implementation

Another paper finds hyperparameter β also helps

VAE Implementation

Another paper finds hyperparameter β also helps

$$\mathcal{L}(\mathbf{X}, \boldsymbol{\theta}) = \frac{m}{n} \sum_{i=1}^n \sum_{j=1}^{d_z} \left(x_{[i],j} - f^{-1} \left(f \left(\mathbf{x}_{[i]}, \boldsymbol{\theta}_e \right), \boldsymbol{\theta}_d \right)_j \right)^2 - \beta \left(\sum_{i=1}^n \sum_{j=1}^{d_z} \mu_{[i],j}^2 + \sigma_{[i],j}^2 - \log(\sigma_{[i],j}^2) - 1 \right)$$

VAE Implementation

```
def L(model, x, m, n, key):
    mu, sigma = model.f(x)
    epsilon = jax.random.normal(key, x.shape[0])
    z = mu + sigma * epsilon
    pred_x = model.f_inverse(z)

    recon = jnp.sum((x - pred_x) ** 2)
    kl = jnp.sum(mu ** 2 + sigma ** 2 - jnp.log(sigma ** 2) -
    1)

    return m / n * recon + kl
```

VAE Implementation

https://colab.research.google.com/drive/1UyR_W6NDIujaJXYlHZh6O3NfaCAMscpH#scrollTo=nmyQ8aE2pSbb

<https://www.youtube.com/watch?v=UZDiGooFs54>

4 Nov - Compensatory rest day (holiday), no lecture

11 Nov - Quiz on autoencoders (not VAE) and Prof. Li on GNNs

18 Nov - Attention and transformers

25 Nov - Reinforcement learning

1 Dec - Foundation models (virtual, optional)