

Autoencoders and Generative Models

CISC 7026: Introduction to Deep Learning

University of Macau

Agenda

1. Review
2. Compression
3. Autoencoders
4. Applications
5. Variational Models
6. Coding

Agenda

1. Review
2. Compression
3. Autoencoders
4. Applications
5. Variational Models
6. Coding

Review

Agenda

1. Review
2. Compression
3. Autoencoders
4. Applications
5. Variational Models
6. Coding

Agenda

1. Review
2. **Compression**
3. Autoencoders
4. Applications
5. Variational Models
6. Coding

Compression

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 concepts with friends (paintings, music, films, etc),
you summarize them

Compression

When you discuss concepts with friends (paintings, music, films, etc), you summarize them

This is a form of **compression**

Compression

When you discuss concepts with friends (paintings, music, films, etc), 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 concepts with friends (paintings, music, films, etc), 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 take some data and reduce its size by removing unnecessary information

Compression

When you discuss concepts with friends (paintings, music, films, etc), 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 take some data and reduce its size by removing unnecessary 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

Fortunately, very talented engineers created the H.264 video codec

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

Fortunately, very talented engineers created the H.264 video codec

H.264 MPEG-AVC **encoder** transforms 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

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

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

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

We download Z from the internet

Compression

We download Z from the internet

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

Compression

We download Z from the internet

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

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

Compression

We download Z from the internet

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

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

Question: How do we watch the video?

Compression

We download Z from the internet

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

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

Question: How do we watch the video?

Answer: Transform or **decode** Z back into pixels

Compression

We download Z from the internet

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

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

Question: How do we watch the video?

Answer: Transform or **decode** Z back into pixels

We need to undo (invert) the encoding function f

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

Compression

We download Z from the internet

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

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

Question: How do we watch the video?

Answer: Transform or **decode** Z back into pixels

We need to undo (invert) the encoding function f

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

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

Compression

To summarize:

Compression

To summarize:

We encode pixels into a bit string to save space

Compression

To summarize:

We encode pixels into a bit string to save space

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

Compression

To summarize:

We encode pixels into a bit string to save space

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

We store the film as a bit string on websites or your computer

Compression

To summarize:

We encode pixels into a bit string to save space

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

We store the film as a bit string on websites or your computer

When you want to watch, we decode the string back into pixels

Compression

To summarize:

We encode pixels into a bit string to save space

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

We store the film as a bit string on websites or your computer

When you want to watch, we decode the string back into pixels

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

Compression

Compression may be **lossy**

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 Models
6. Coding

Agenda

1. Review
2. Compression
3. **Autoencoders**
4. Applications
5. Variational Models
6. Coding

Autoencoders

The encoders and decoders for images, videos, etc are very complex

Autoencoders

The encoders and decoders for images, videos, etc are very complex

Neural networks can represent any continuous function

Autoencoders

The encoders and decoders for images, videos, etc are very complex

Neural networks can represent any continuous function

Let us learn neural network encoders and decoders

Autoencoders

The encoders and decoders for images, videos, etc are very complex

Neural networks can represent any continuous function

Let us learn neural network encoders and decoders

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

Autoencoders

The encoders and decoders for images, videos, etc are very complex

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

The encoders and decoders for images, videos, etc are very complex

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

The encoders and decoders for images, videos, etc are very complex

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**

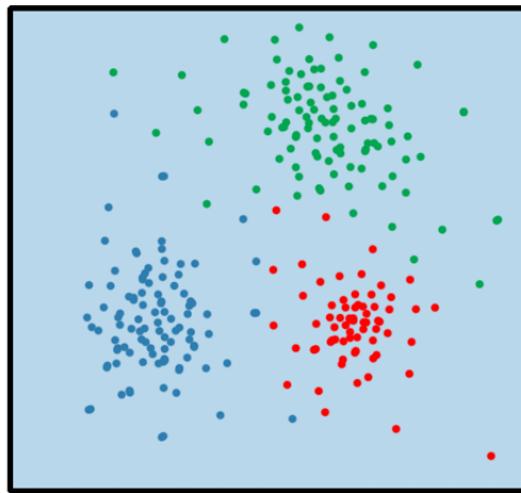
Notice there is no Y this time

Training autoencoders is different than what we have seen before

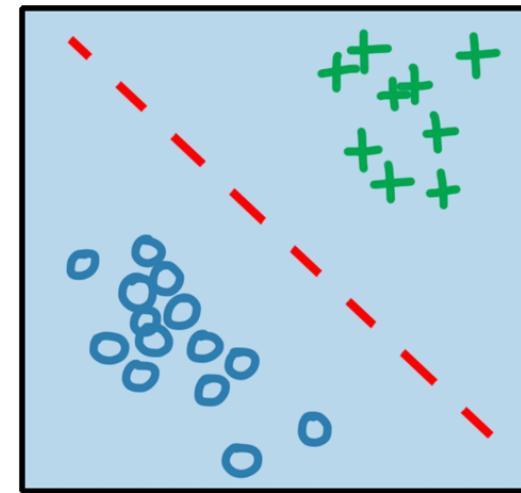
Autoencoders

machine learning

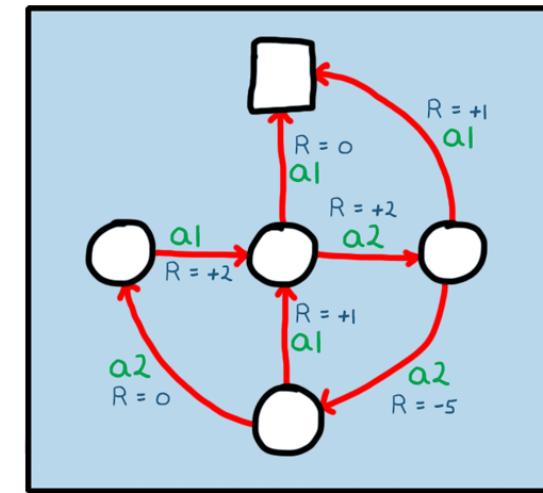
unsupervised
learning



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**

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

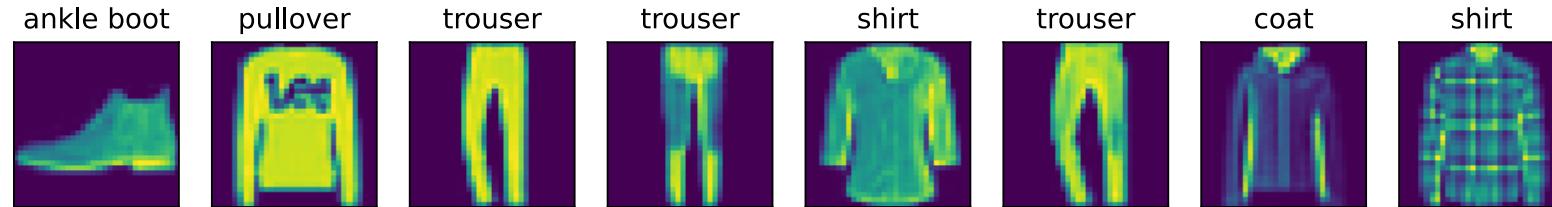
The training algorithm may generate labels

Autoencoders

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

Autoencoders

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



$$d_x : 28 \times 28$$

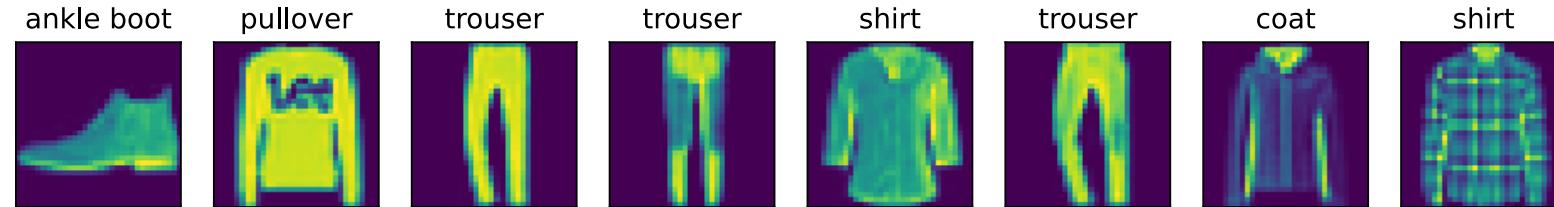
$$d_z : 4$$

$$X : [0, 1]^{d_x}$$

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

Autoencoders

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



$$d_x : 28 \times 28$$

$$d_z : 4$$

$$X : [0, 1]^{d_x}$$

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

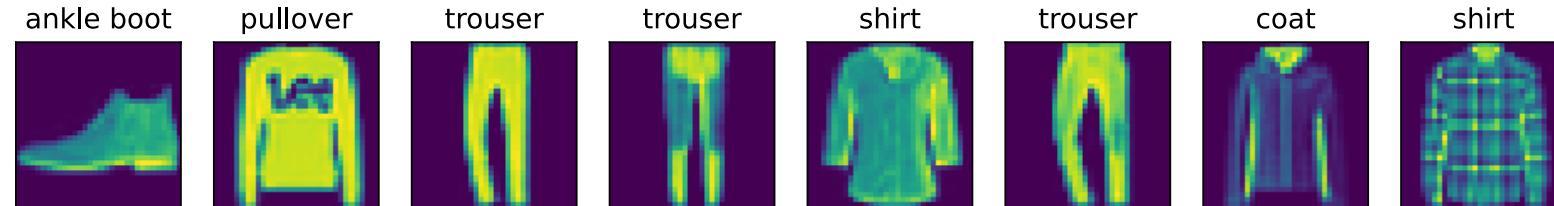
$$f(x, \theta) = z$$

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

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

Autoencoders

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



$$d_x : 28 \times 28$$

$$d_z : 4$$

$$X : [0, 1]^{d_x}$$

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

$$f(x, \theta) = z$$

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

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

How do we find θ ?

Autoencoders

Let us find f , then find the inverse

Autoencoders

Let us find f , then find the inverse

Start with a perceptron

Autoencoders

Let us find f , then find the inverse

Start with a perceptron

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

Autoencoders

Let us find f , then find the inverse

Start with a perceptron

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

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

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

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

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

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

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

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

Autoencoders

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

Autoencoders

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

Question: Any issues?

Hint: What if $d_x \neq d_z$?

Autoencoders

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

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

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

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

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

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

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

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

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

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

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

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

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

Autoencoders

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

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

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

Autoencoders

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

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

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

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

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

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 Models
6. Coding

Agenda

1. Review
2. Compression
3. Autoencoders
4. **Applications**
5. Variational Models
6. Coding

Applications

We can use autoencoders for other tasks too

Applications

We can use autoencoders for other tasks too

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

Applications

We can use autoencoders for other tasks too

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



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

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

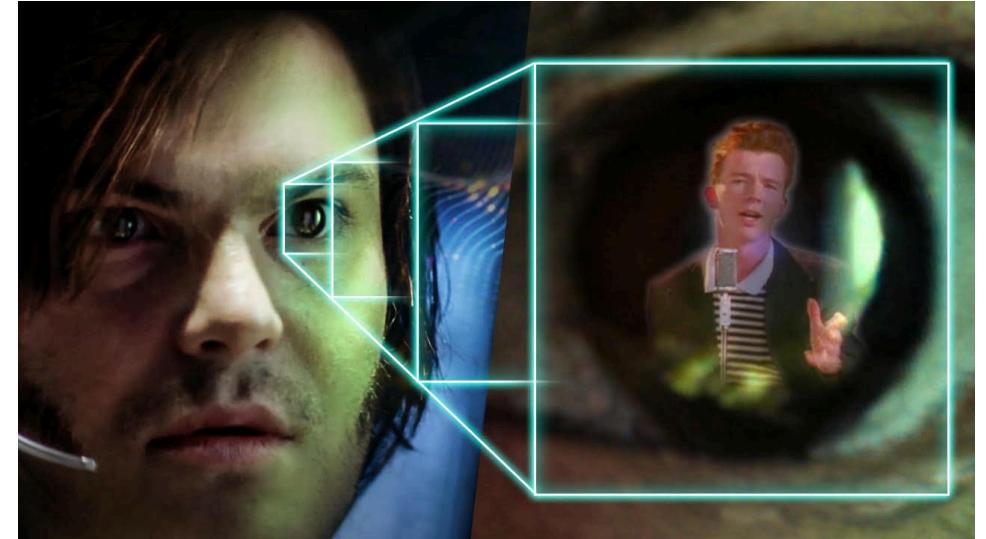
$$\mathcal{L}(\mathbf{X}, \boldsymbol{\theta}) = \sum_{i=1}^n \sum_{j=1}^{d_x} \left(x_{[i],j} - f^{-1} \left(f \left(\text{blur} \left(\mathbf{x}_{[i]} + \varepsilon \right), \boldsymbol{\theta}_e \right), \boldsymbol{\theta}_d \right)_j \right)^2$$

Applications

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

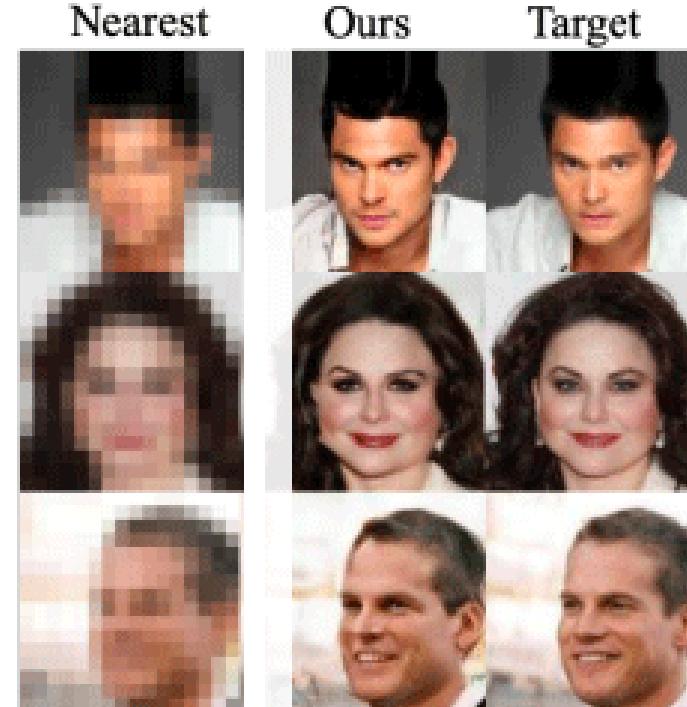
Applications

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

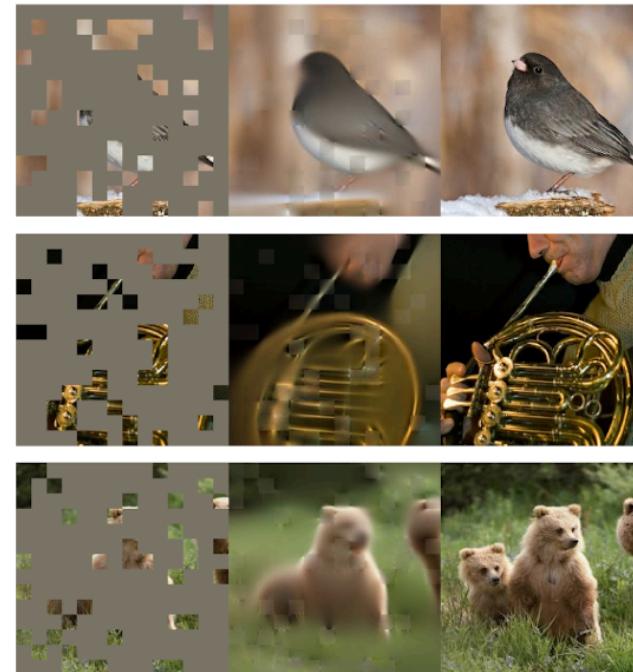
We can even hide parts of the image

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

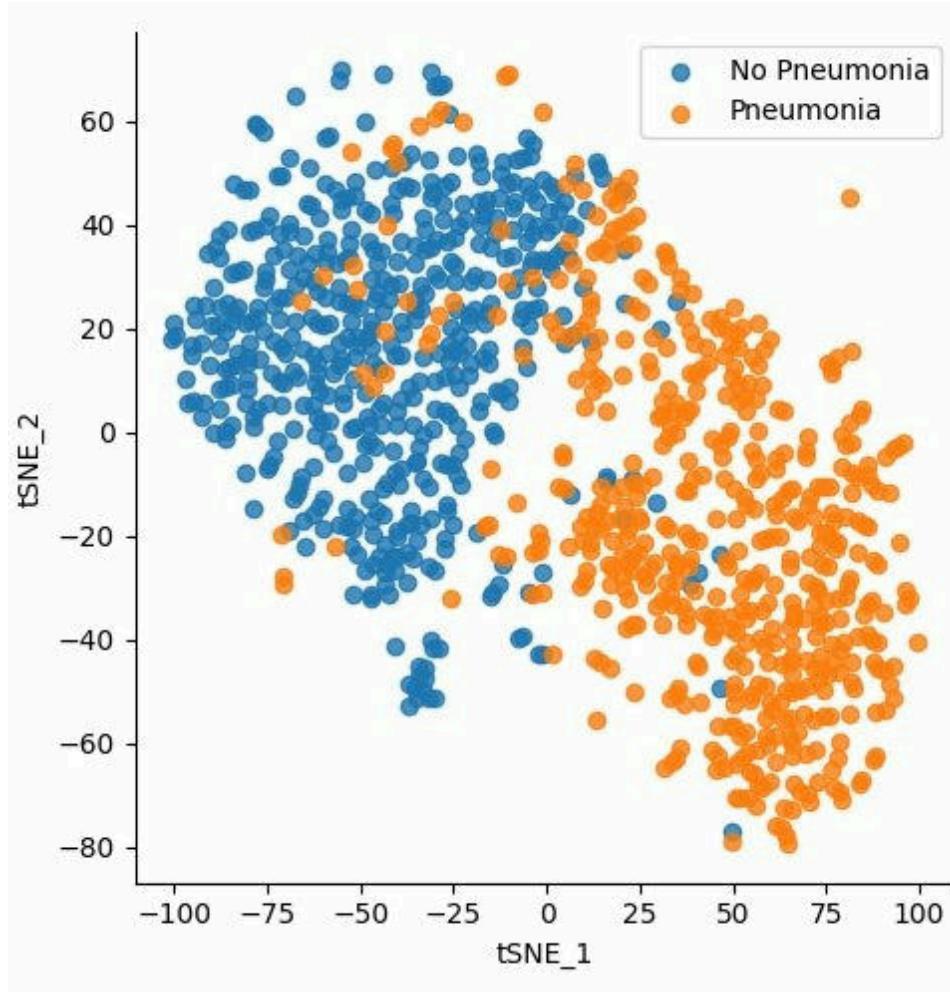
What is happening here? How can these models do this?

Applications

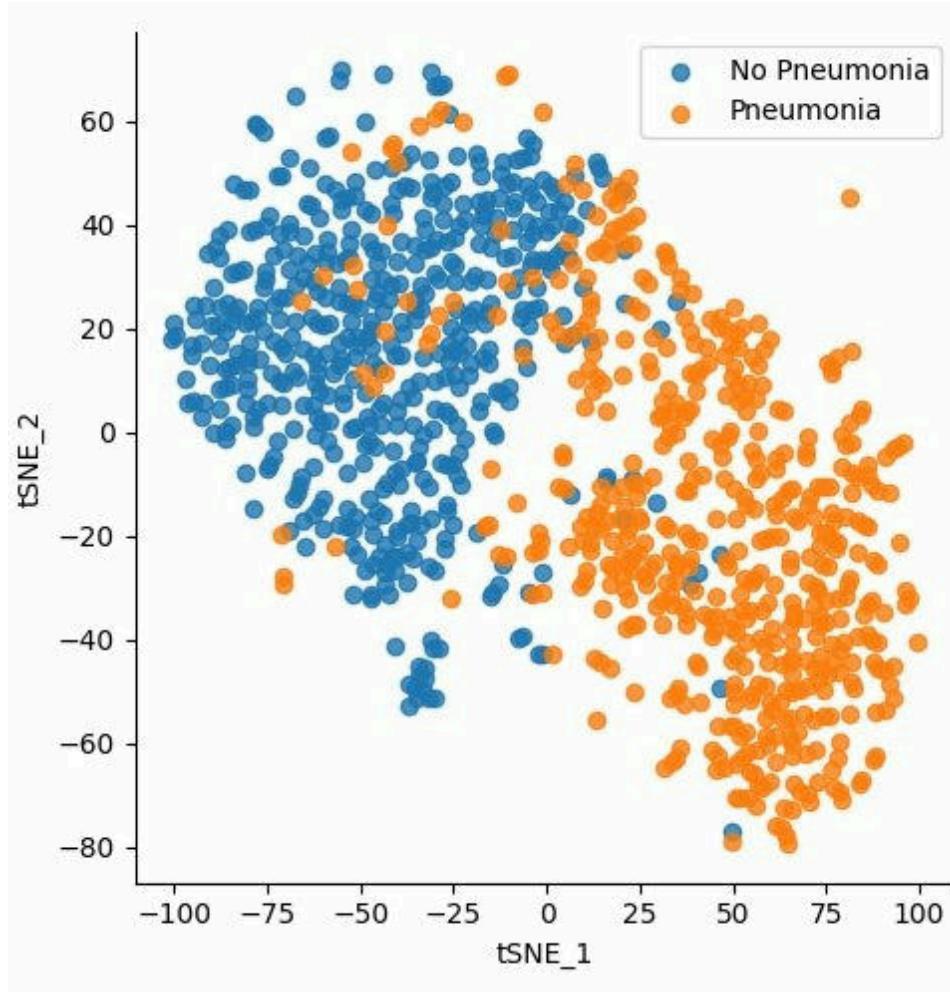
What is happening here? How can these models do this?

Autoencoders learn the structure of the dataset

Applications

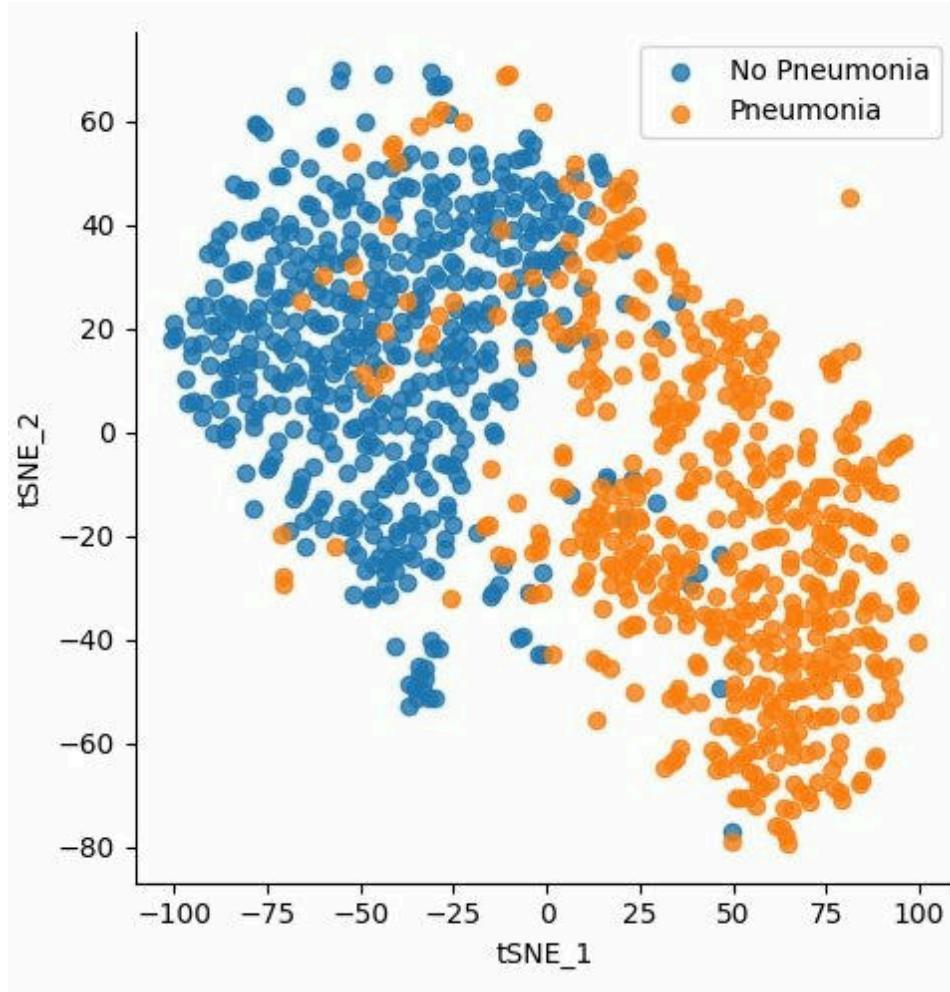


Applications



X : Pictures of human lungs

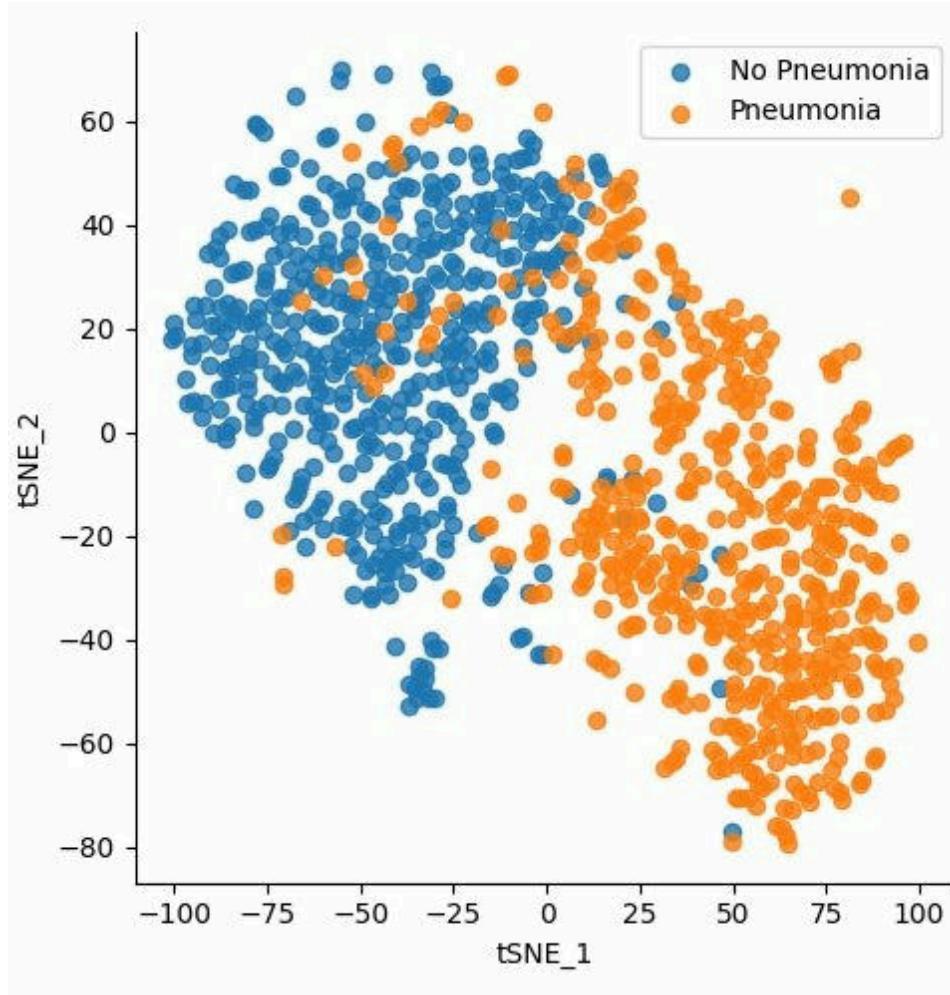
Applications



X : Pictures of human lungs

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

Applications

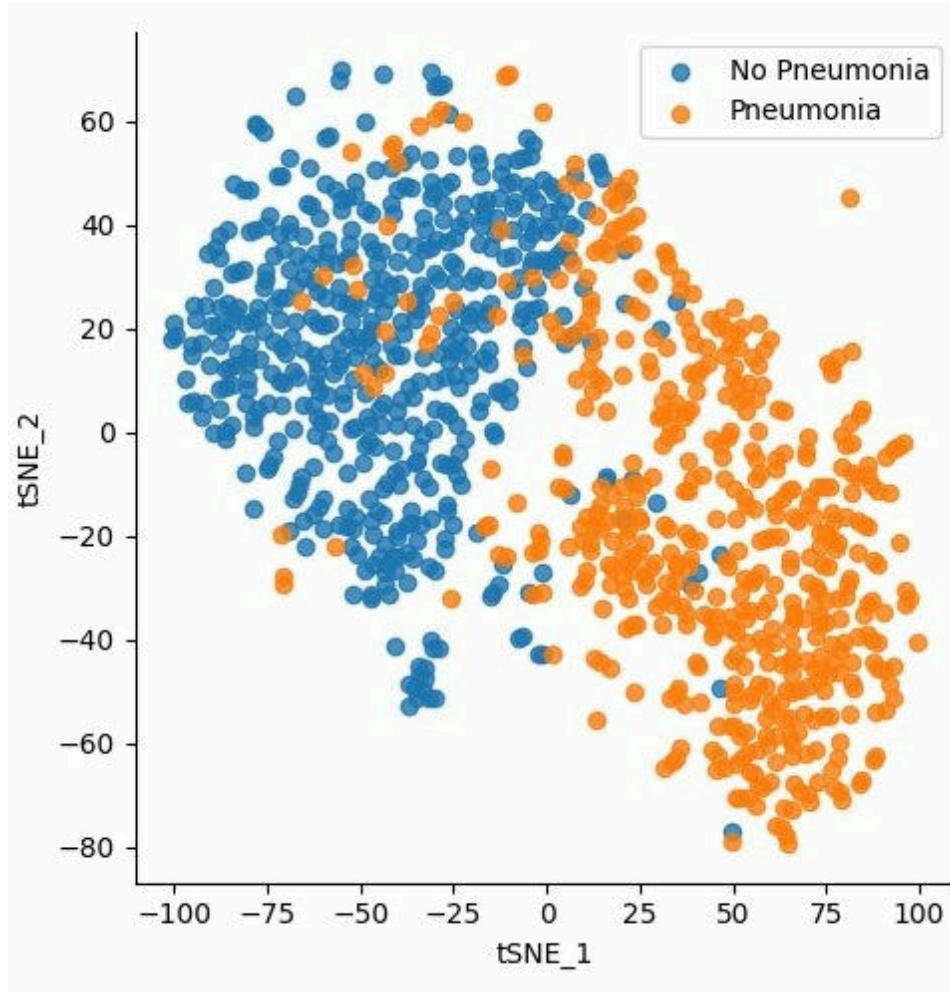


X : Pictures of human lungs

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

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 pictures from our world, then the autoencoders learn the structure of the world

Applications

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

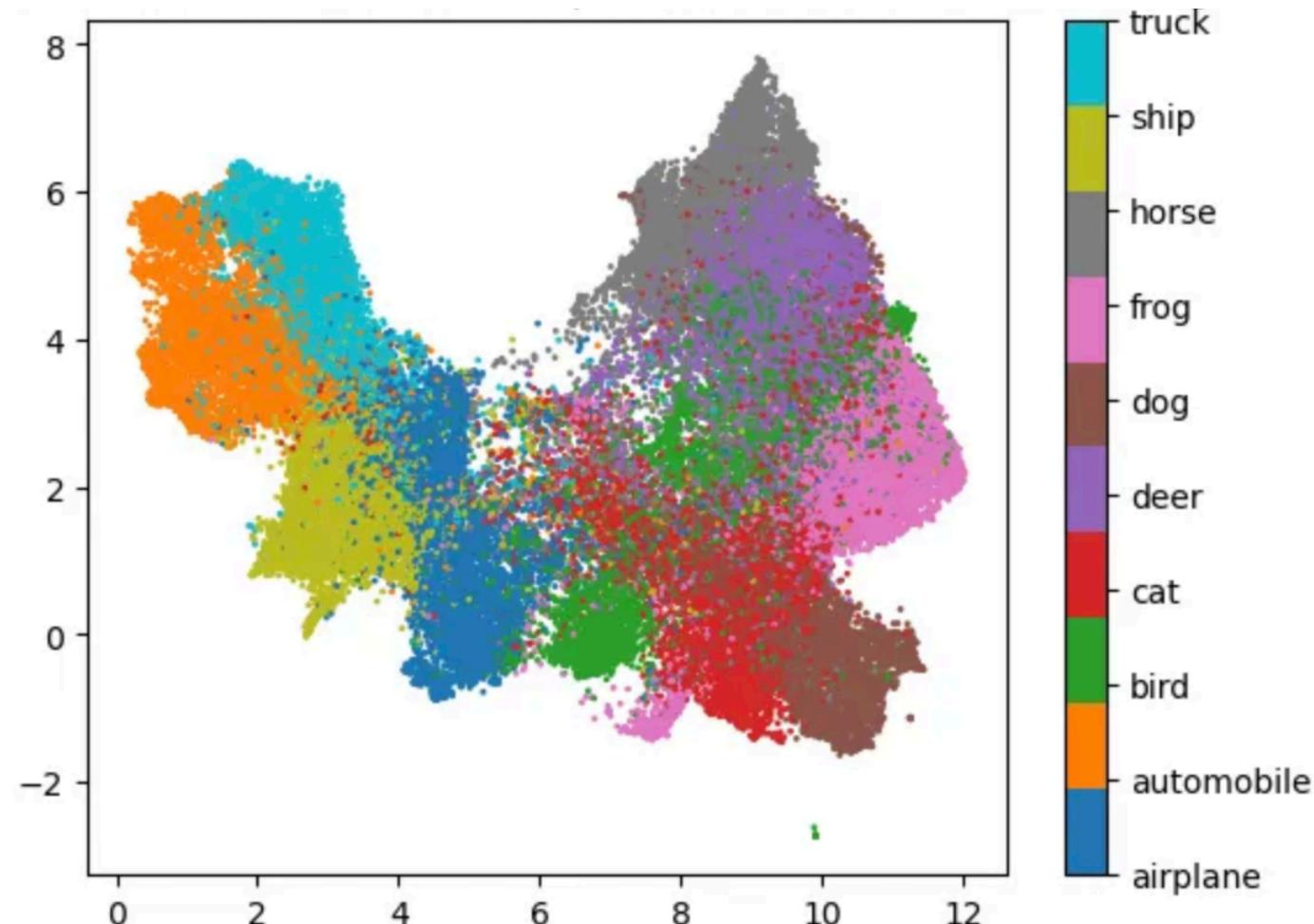
Applications

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

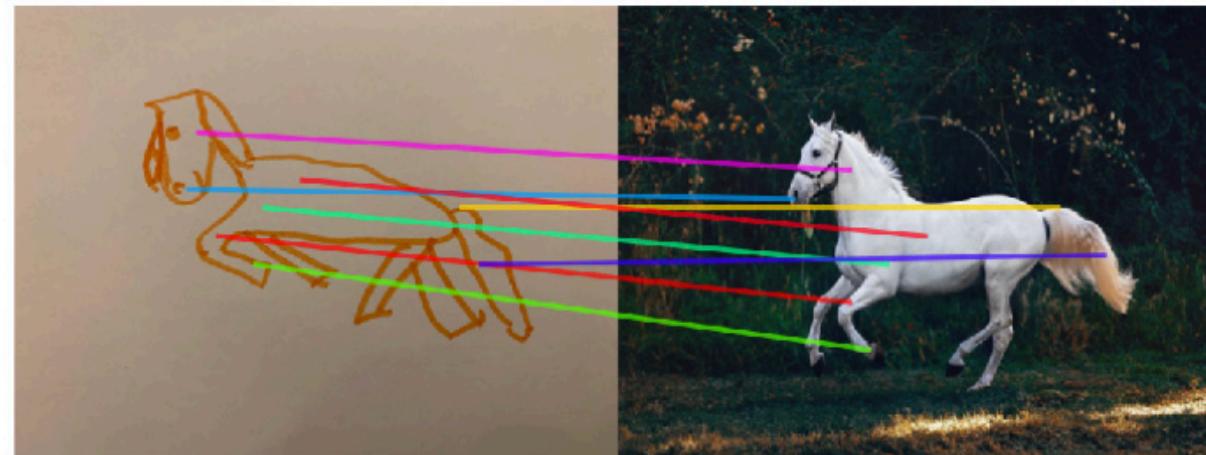
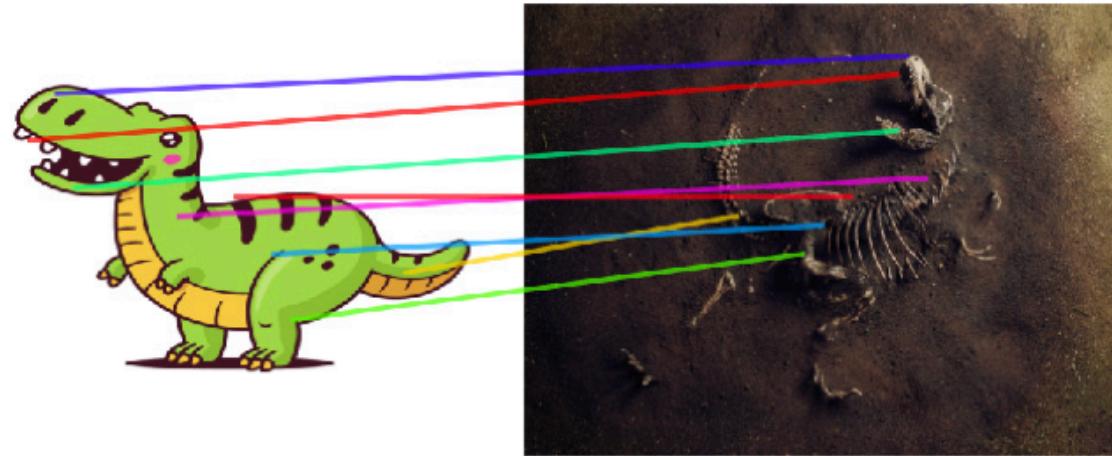
Nobody tells them what a dog or cat is

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

Clearly, these networks can understand our world

Agenda

1. Review
2. Compression
3. Autoencoders
4. **Applications**
5. Variational Models
6. Coding

Agenda

1. Review
2. Compression
3. Autoencoders
4. Applications
5. **Variational Models**
6. Coding

Variational Models



Variational Models



These pictures were created by an autoencoder

Variational Models



These pictures were created by an autoencoder

But these people do not exist!

Variational Models

Autoencoders are useful for compression and denoising

Variational Models

Autoencoders are useful for compression and denoising

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

Variational Models

Autoencoders are useful for compression and denoising

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

A generative model is a “creative” model that creates new data

Variational Models

Autoencoders are useful for compression and denoising

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

A generative model is a “creative” model that creates new data

- Generate pictures of people that do not exist

Variational Models

Autoencoders are useful for compression and denoising

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

A generative model is a “creative” model that creates new data

- Generate pictures of people that do not exist
- Write a new book

Variational Models

Autoencoders are useful for compression and denoising

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

A generative model is a “creative” model that creates new data

- Generate pictures of people that do not exist
- Write a new book
- Create new proteins

Variational Models

Autoencoders are useful for compression and denoising

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

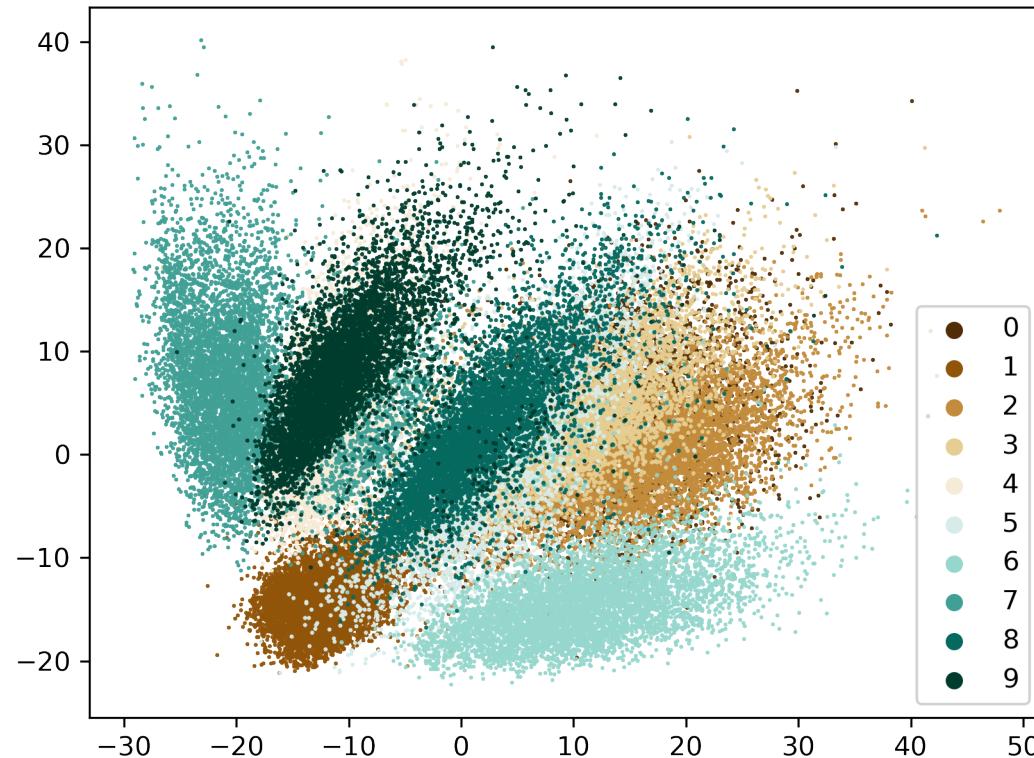
A generative model is a “creative” model that creates new data

- Generate pictures of people that do not exist
- Write a new book
- Create new proteins

How does this work?

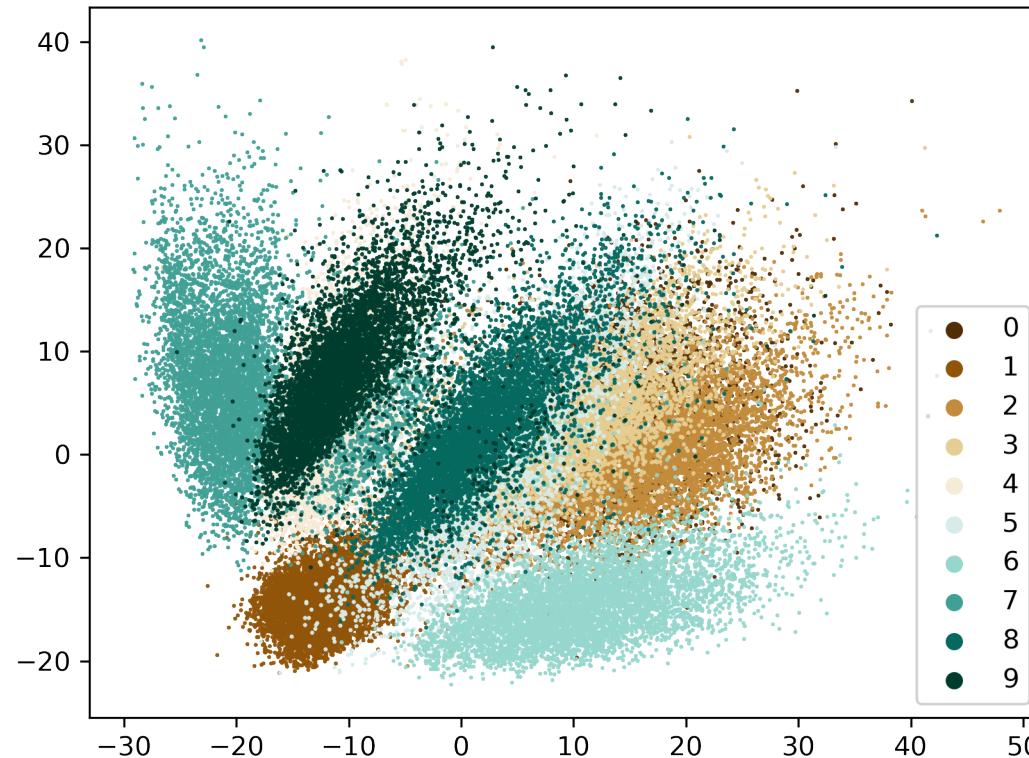
Variational Models

Latent space Z after training on the clothes dataset with $d_z = 2$



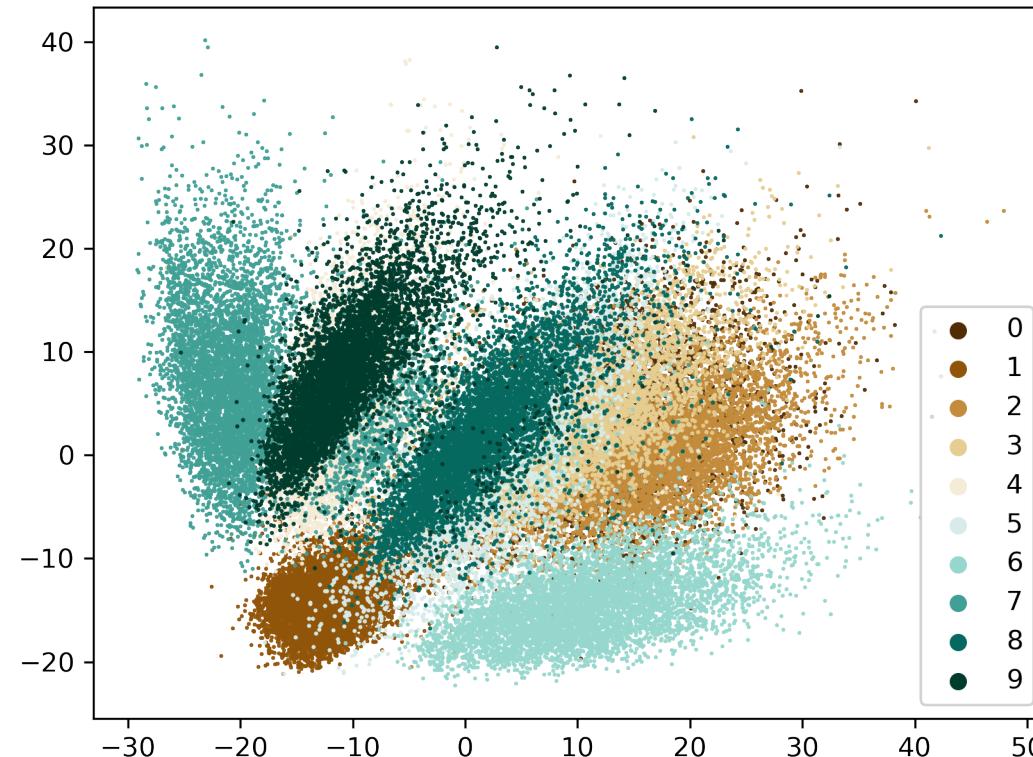
Variational Models

Latent space Z after training on the clothes dataset with $d_z = 2$

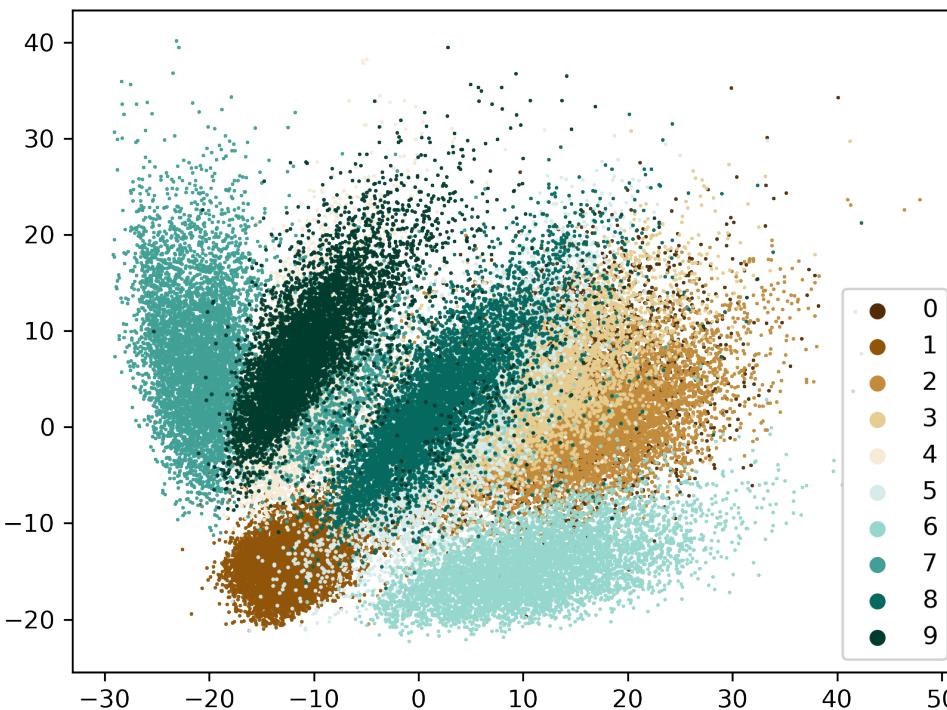


Variational Models

What happens if we decode a new point?

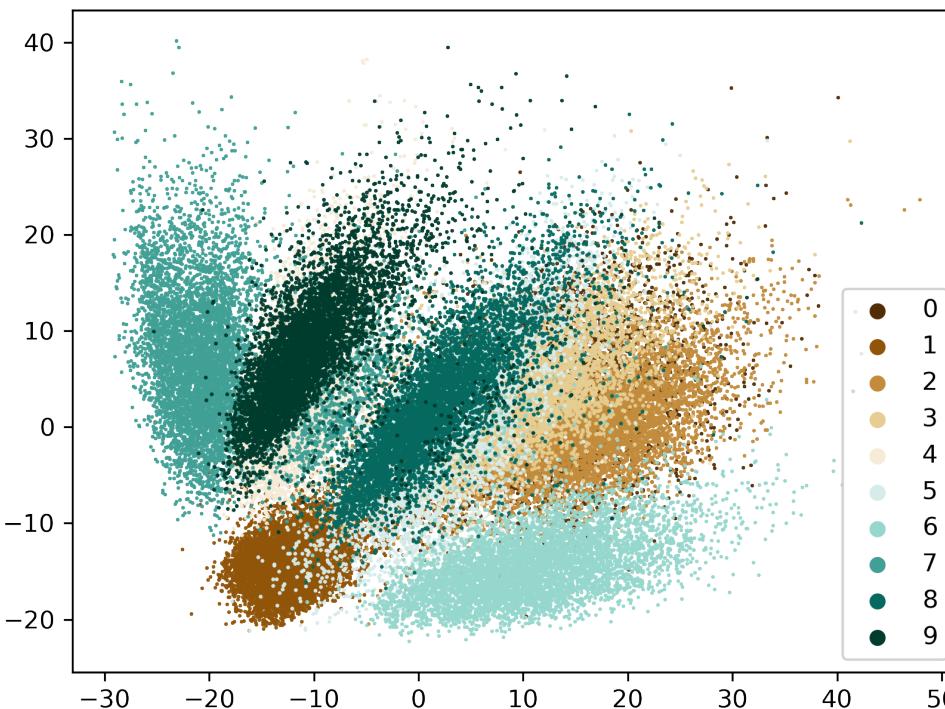


Variational Models



Autoencoder generative model:

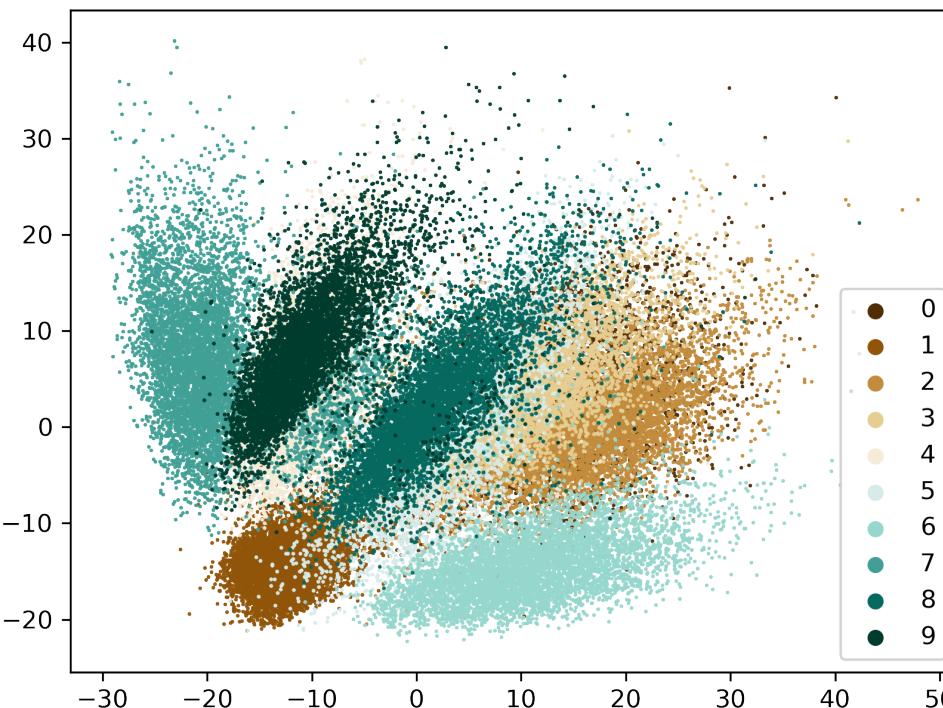
Variational Models



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 Models

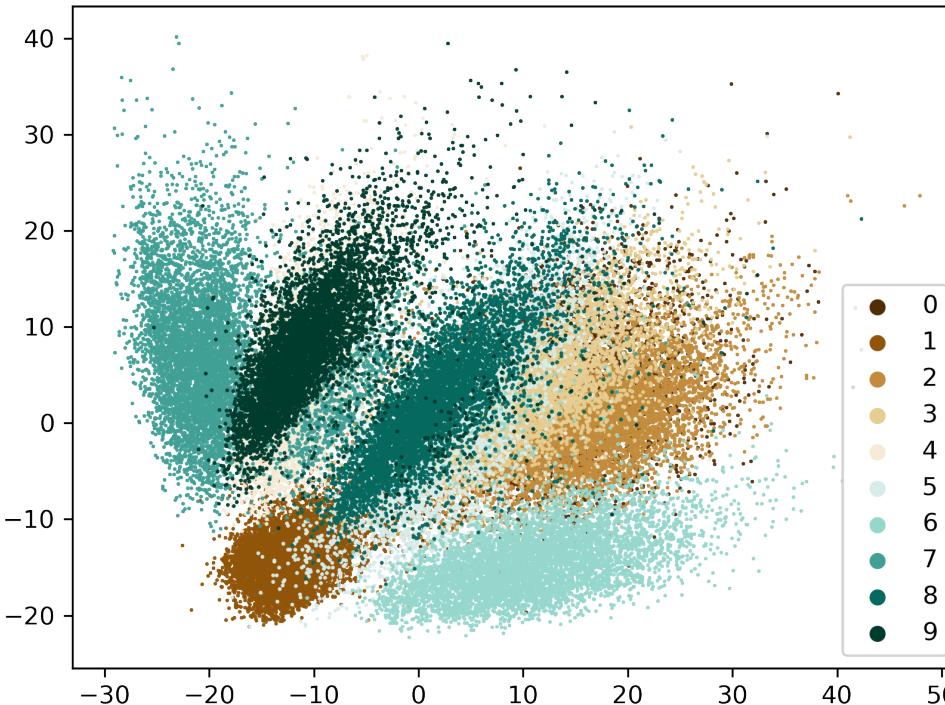


Autoencoder generative model:

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 Models



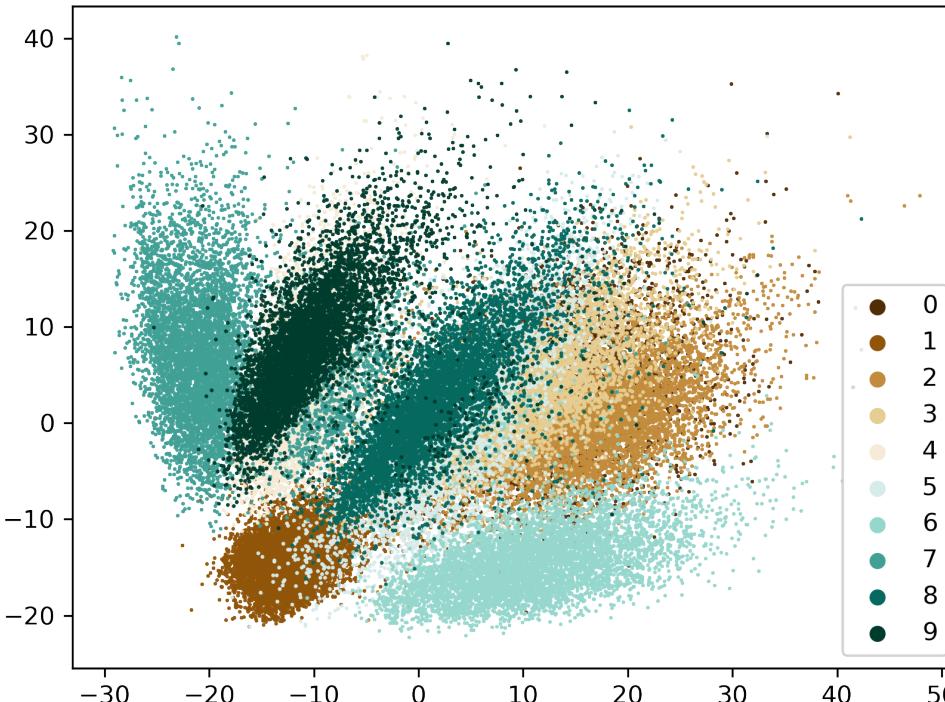
Autoencoder generative model:

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 Models



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

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

Variational Models



Variational Models



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

But there is a problem

But there is a problem

As d_z increases, it becomes
difficult to find useful parts of
latent space

But there is a problem

As d_z increases, it becomes
difficult to find useful parts of
latent space

For large d_z , similar inputs map to
 z that are very far from each other

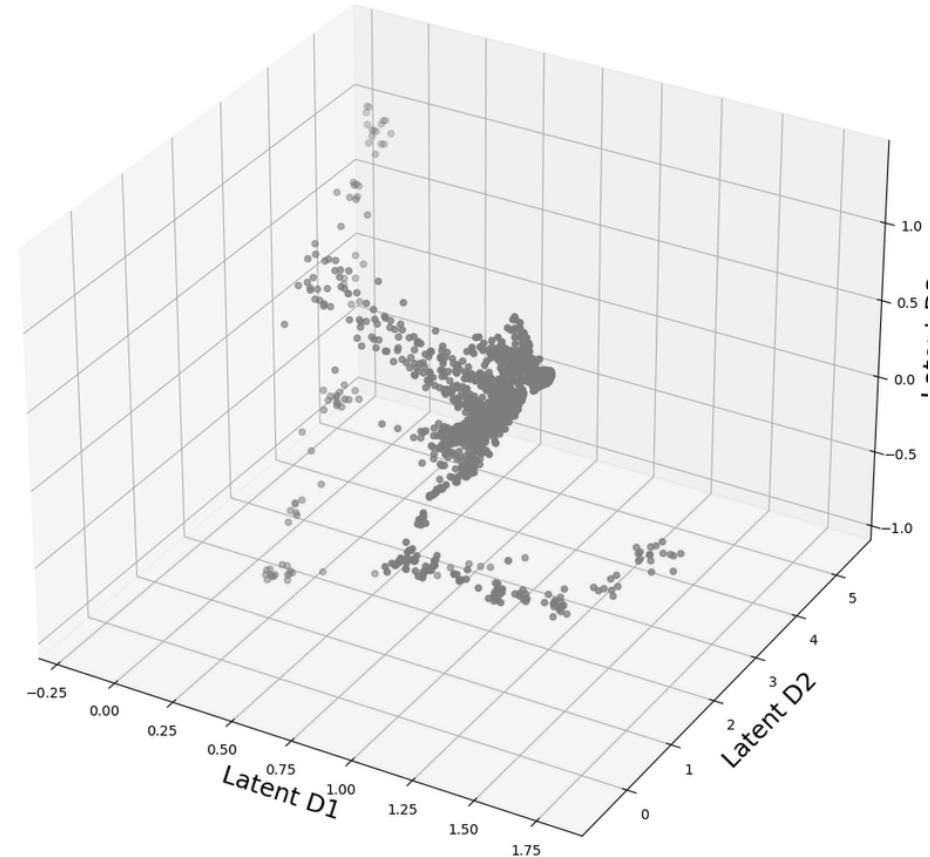
Variational Models

But there is a problem

As d_z increases, it becomes difficult to find useful parts of latent space

For large d_z , similar inputs map to z that are very far from each other

Points $z + \varepsilon$ decode into garbage



Question: What can we do?

Question: What can we do?

Answer: Force the points to be close together!

Variational Models

Question: What can we do?

Answer: Force the points to be close together!

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

Variational Models

VAE discovered by Diederik Kingma (also adam optimizer)

Variational Models

VAE discovered by Diederik Kingma (also adam optimizer)



Variational Models

Variational autoencoders (VAEs) do three things:

Variational Models

Variational autoencoders (VAEs) do three things:

1. Make it easy to sample random z

Variational Models

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 Models

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 Models

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 Models

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

Variational Models

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

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

Variational Models

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

Variational Models

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

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

Variational Models

If $z_{[i]}, \dots, z_{[n]}$ 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 Models

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

Variational Models

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

We have to remember conditional probabilities

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

Variational Models

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

We have to remember conditional probabilities

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

Variational Models

Introduce one more Bayesian concept called **marginalization**

Variational Models

Introduce one more Bayesian concept called **marginalization**

Assume we have the probability of two events

$$P(A \cap B) = P(A, B)$$

Variational Models

Introduce one more Bayesian concept called **marginalization**

Assume we have the probability of two events

$$P(A \cap B) = P(A, B)$$

From Bayes rule

$$P(A \cap B) = P(A \mid B)P(B)$$

Variational Models

Introduce one more Bayesian concept called **marginalization**

Assume we have the probability of two events

$$P(A \cap B) = P(A, B)$$

From Bayes rule

$$P(A \cap B) = P(A \mid B)P(B)$$

We can find $P(A)$ by marginalizing out B

Variational Models

Introduce one more Bayesian concept called **marginalization**

Assume we have the probability of two events

$$P(A \cap B) = P(A, B)$$

From Bayes rule

$$P(A \cap B) = P(A \mid B)P(B)$$

We can find $P(A)$ by marginalizing out B

$$P(A) = \int_B P(A \mid B)P(B) \mathrm{d}B$$

Variational Models

Warning: Derivation of variational autoencoders is difficult

Variational Models

Warning: Derivation of variational autoencoders is difficult

I have done my best to simplify it, but it is still complex

Variational Models

Warning: Derivation of variational autoencoders is difficult

I have done my best to simplify it, but it is still complex

There might be small mistakes!

Key idea:

Step 1: Encoder produces a distribution $P(z)$ from an input x

$$P(z \mid x)$$

Variational Models

Key idea:

Step 1: Encoder produces a distribution $P(z)$ from an input x

$$P(z \mid x)$$

Step 2: Decoder produces x_{new} sampled from a distribution

$$x_{\text{new}} \sim P(x \mid z)P(z)$$

Variational Models

Key idea:

Step 1: Encoder produces a distribution $P(z)$ from an input x

$$P(z \mid x)$$

Step 2: Decoder produces x_{new} sampled from a distribution

$$x_{\text{new}} \sim P(x \mid z)P(z)$$

VAE is generative because we sample from a distribution

Variational Models

Key idea:

Step 1: Encoder produces a distribution $P(z)$ from an input x

$$P(z \mid x)$$

Step 2: Decoder produces x_{new} sampled from a distribution

$$x_{\text{new}} \sim P(x \mid z)P(z)$$

VAE is generative because we sample from a distribution

We are sampling new z that were not in the dataset!

Variational Models

Let us start with the encoder

Variational Models

Let us start with the encoder

We want to turn an input x into a latent normal distribution over Z

Variational Models

Let us start with the encoder

We want to turn an input x into a latent normal distribution over Z

Start with the joint probability $P(z, x)$

$$P(z, x) = \underbrace{P(z | x)}_{\text{Likelihood}} \underbrace{P(x)}_{\text{Prior}}$$

$$P(x)$$

Probability of a specific input
(image) existing in the world

Variational Models

Let us start with the encoder

We want to turn an input x into a latent normal distribution over Z

Start with the joint probability $P(z, x)$

$$P(z, x) = \underbrace{P(z | x)}_{\text{Likelihood}} \underbrace{P(x)}_{\text{Prior}}$$

$P(x)$

Probability of a specific input
(image) existing in the world

$P(z | x)$

Latent distribution that
summarizes x

Variational Models

We can marginalize to get the latent distribution $P(z)$

$$P(z) = \int_x P(z \mid x) P(x) dx$$

We approximate this function using the encoder

$$f(x, \theta_e) = P(z)$$

Question: What distribution should $P(z)$ be?

Variational Models

We can marginalize to get the latent distribution $P(z)$

$$P(z) = \int_x P(z \mid x) P(x) dx$$

We approximate this function using the encoder

$$f(x, \theta_e) = P(z)$$

Question: What distribution should $P(z)$ be?

Answer: A normal distribution $\mathcal{N}(\mu, \sigma)$

Variational Models

We represent a normal distribution with a **mean** $\mu \in M$ and a **standard deviation** $\sigma \in \Sigma$

Variational Models

We represent a normal distribution with a **mean** $\mu \in M$ and a **standard deviation** $\sigma \in \Sigma$

$$f : X \times \Theta \mapsto M \times \Sigma$$

Variational Models

We represent a normal distribution with a **mean** $\mu \in M$ and a **standard deviation** $\sigma \in \Sigma$

$$f : X \times \Theta \mapsto M \times \Sigma$$

So our encoder should output two values

$$M \in \mathbb{R}^{d_z}, \Sigma \in \mathbb{R}_+^{d_z}$$

Variational Models

```
core = nn.Sequential(...)  
mu_layer = nn.Linear(d_h, d_z)  
# Neural networks output real numbers  
# But sigma must be positive  
# So we output log sigma, because e^(sigma) is always  
positive  
log_sigma_layer = nn.Linear(d_h, d_z)  
  
tmp = core(x)  
mu = mu_layer(tmp)  
log_sigma = log_sigma_layer(tmp)  
dist = (mu, exp(sigma))
```

Variational Models

Decoder maps the distribution $P(z)$ to $P(x)$

$$P(x) = \int_z P(x \mid z)P(z)$$

$$x \sim P(x)$$

$$P(z)$$

Normal dist. given by encoder

Variational Models

Decoder maps the distribution $P(z)$ to $P(x)$

$$P(x) = \int_z P(x \mid z)P(z)$$

$$x \sim P(x)$$

$P(z)$ Normal dist. given by encoder

$P(x \mid z)$ Output distribution

x Output (e.g., an image)

Variational Models

$$P(\mathbf{x}) = \int_z P(\mathbf{x} \mid \mathbf{z})P(\mathbf{z})$$

$$\mathbf{x} \sim P(\mathbf{x})$$

Integral is intractable, so the decoder inputs and outputs a **sample**

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

Variational Models

```
decoder = nn.Sequential(  
    nn.Linear(d_z, d_h),  
    ...  
    nn.Linear(d_h, d_x)  
)
```

Decoder is the same as a normal autoencoder

Variational Models

So far, we have the encoder

$$f : X \times \Theta \mapsto M \times \Sigma$$

Variational Models

So far, we have the encoder

$$f : X \times \Theta \mapsto M \times \Sigma$$

And the decoder

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

Variational Models

So far, we have the encoder

$$f : X \times \Theta \mapsto M \times \Sigma$$

And the decoder

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

Question: Any problems?

Variational Models

So far, we have the encoder

$$f : X \times \Theta \mapsto M \times \Sigma$$

And the decoder

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

Question: Any problems?

Answer: The encoder output is a distribution $M \times \Sigma$, but the decoder input is a sample Z

Variational Models

To convert from a distribution over Z to a z , we must sample

Variational Models

To convert from a distribution over Z to a z , we must sample

$$z \sim P(z)$$

Variational Models

To convert from a distribution over Z to a z , we must sample

$$z \sim P(z)$$

But sampling from a distribution is not differentiable

Variational Models

To convert from a distribution over Z to a z , we must sample

$$z \sim P(z)$$

But sampling from a distribution is not differentiable

The authors discover the **reparameterization trick**

Variational Models

To convert from a distribution over Z to a z , we must sample

$$z \sim P(z)$$

But sampling from a distribution is not differentiable

The authors discover the **reparameterization trick**

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

Variational Models

To convert from a distribution over Z to a z , we must sample

$$z \sim P(z)$$

But sampling from a distribution is not differentiable

The authors discover the **reparameterization trick**

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

Gradient can flow through μ, σ

Variational Models

To convert from a distribution over Z to a z , we must sample

$$z \sim P(z)$$

But sampling from a distribution is not differentiable

The authors discover the **reparameterization trick**

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

Gradient can flow through μ, σ

We can draw samples but still use gradient descent

Put it all together

Variational Models

Put it all together

Step 1: Encode the input to a latent distribution

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

Step 2: Generate a sample from distribution

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

Step 3: Decode the sample

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

Variational Models

Todo loss

Variational Models

We want to constrain the distribution

$$P(z \mid x) \approx \mathcal{N}(0, 1)$$

Variational Models

We want to constrain the distribution

$$P(z \mid x) \approx \mathcal{N}(0, 1)$$

We will implement this constraint in the loss function

Variational Models

We want to constrain the distribution

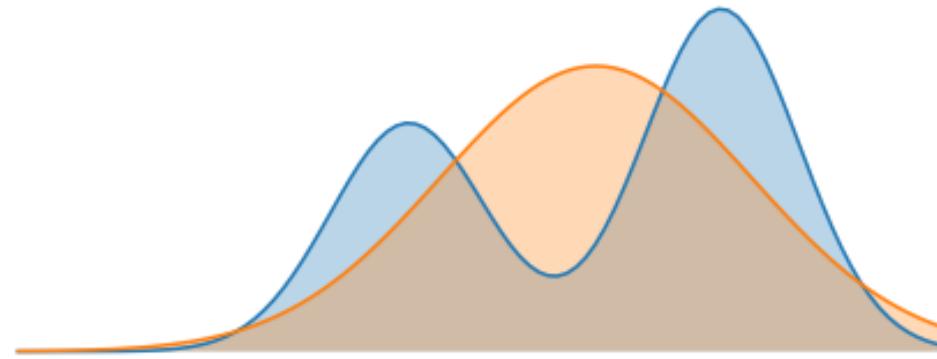
$$P(z \mid x) \approx \mathcal{N}(0, 1)$$

We will implement this constraint in the loss function

Question: How to measure distance between two distributions?

Variational Models

Answer: KL divergence



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

Variational Models

$$\mathcal{L}_{\text{KL}}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^{d_z} \text{KL}\left(P(z_j \mid \mathbf{x}), \mathcal{N}(\mathbf{0}, \mathbf{1})\right)$$

$$\mathcal{L}_{\text{KL}}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^{d_z} \text{KL}\left(f(\mathbf{x}, \boldsymbol{\theta}_e)_j, \mathcal{N}(\mathbf{0}, \mathbf{1})\right)$$

This ensures that the latent representations \mathbf{z} are easy to sample

Variational Models

$$\mathcal{L}_{\text{KL}}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^{d_z} \text{KL}\left(P(z_j \mid \mathbf{x}), \mathcal{N}(0, 1)\right)$$

Remember, our neural network outputs $P(z)$ as μ, σ

$$f(\mathbf{x}, \boldsymbol{\theta}) = \mu, \sigma$$

Variational Models

$$\mathcal{L}_{\text{KL}}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^{d_z} \text{KL}\left(P(z_j \mid \mathbf{x}), \mathcal{N}(0, 1)\right)$$

Remember, our neural network outputs $P(\mathbf{z})$ as μ, σ

$$f(\mathbf{x}, \boldsymbol{\theta}) = \mu, \sigma$$

The KL divergence between two normal distributions simplifies to

$$\mathcal{L}_{\text{KL}}(\mathbf{x}, \boldsymbol{\theta}) = \left(\sum_{j=1}^{d_z} 1 + \sigma_j^2 + \mu_j^2 - \exp(\sigma_j^2) \right)$$

Variational Models

$$\mathcal{L}_{\text{KL}}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^{d_z} \text{KL}\left(P(z_j \mid \mathbf{x}), \mathcal{N}(0, 1)\right)$$

Remember, our neural network outputs $P(\mathbf{z})$ as μ, σ

$$f(\mathbf{x}, \boldsymbol{\theta}) = \mu, \sigma$$

The KL divergence between two normal distributions simplifies to

$$\mathcal{L}_{\text{KL}}(\mathbf{x}, \boldsymbol{\theta}) = \left(\sum_{j=1}^{d_z} 1 + \sigma_j^2 + \mu_j^2 - \exp(\sigma_j^2) \right)$$

Variational Models

$$\mathcal{L}_{\text{KL}}(x, \theta) = \left(\sum_{j=1}^{d_z} 1 + \sigma_j^2 + \mu_j^2 - \exp(\sigma_j^2) \right)$$

Variational Models

$$\mathcal{L}_{\text{KL}}(\boldsymbol{x}, \boldsymbol{\theta}) = \left(\sum_{j=1}^{d_z} 1 + \sigma_j^2 + \mu_j^2 - \exp(\sigma_j^2) \right)$$

Over the whole dataset

$$\mathcal{L}_{\text{KL}}(\boldsymbol{X}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(\sum_{j=1}^{d_z} 1 + \sigma_{[i],j}^2 + \mu_{[i],j}^2 - \exp(\sigma_{[i],j}^2) \right)$$

Variational Models

$$\mathcal{L}_{\text{KL}}(\mathbf{X}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(\sum_{j=1}^{d_z} 1 + \sigma_{[i],j}^2 + \mu_{[i],j}^2 - \exp(\sigma_{[i],j}^2) \right)$$

Variational Models

$$\mathcal{L}_{\text{KL}}(\mathbf{X}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(\sum_{j=1}^{d_z} 1 + \sigma_{[i],j}^2 + \mu_{[i],j}^2 - \exp(\sigma_{[i],j}^2) \right)$$

The KL loss ensures that $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$

Variational Models

$$\mathcal{L}_{\text{KL}}(\mathbf{X}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(\sum_{j=1}^{d_z} 1 + \sigma_{[i],j}^2 + \mu_{[i],j}^2 - \exp(\sigma_{[i],j}^2) \right)$$

The KL loss ensures that $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$

We also want to make sure we reconstruct the input!

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

Variational Models

$$\mathcal{L}_{\text{KL}}(\mathbf{X}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(\sum_{j=1}^{d_z} 1 + \sigma_{[i],j}^2 + \mu_{[i],j}^2 - \exp(\sigma_{[i],j}^2) \right)$$

The KL loss ensures that $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$

We also want to make sure we reconstruct the input!

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

cal(L)(bold(X), bold(theta))

Variational Models

The loss for a VAE contains the reconstruction error and the KL loss