

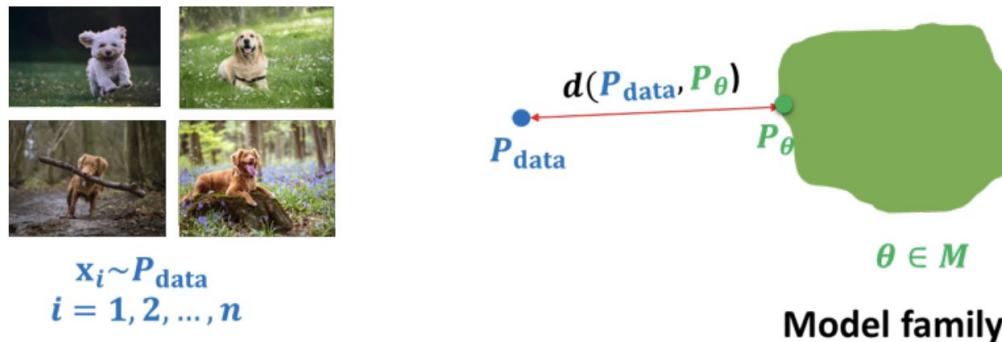
Introduction to Data Mining

CS 145

Lecture 10:

Self-Supervised Learning: (Variational)
AutoEncoder and Contrastive Learning

- We are given a training set of examples, e.g., images of dogs



- We want to learn a probability distribution $p(x)$ over images x such that
 - Generation:** If we sample $x_{\text{new}} \sim p(x)$, x_{new} should look like a dog (*sampling*)
 - Density estimation:** $p(x)$ should be high if x looks like a dog, and low otherwise (*anomaly detection*)
 - Unsupervised representation learning:** We should be able to learn what these images have in common, e.g., ears, tail, etc. (*features*)
- First question: how to represent $p_\theta(x)$. Second question: **how to learn it.**

$$\log p(\mathcal{D}; \theta) = \sum_{n=1}^N \log p(\mathbf{x}^{(n)}; \theta) = \sum_{n=1}^N \log \left(\sum_{z^{(n)}} p(z^{(n)}, \mathbf{x}^{(n)}; \theta) \right)$$

- We introduce auxilliary distributions $q_n(z^{(n)})$ over each of the latent variables

$$\begin{aligned} \sum_{n=1}^N \log \left(\sum_{z^{(n)}} p(z^{(n)}, \mathbf{x}^{(n)}; \theta) \right) &= \sum_{n=1}^N \log \left(\sum_{z^{(n)}} q_n(z^{(n)}) \frac{p(z^{(n)}, \mathbf{x}^{(n)}; \theta)}{q_n(z^{(n)})} \right) \\ &= \sum_{n=1}^N \log \left(\mathbb{E}_{q_n(z^{(n)})} \left[\frac{p(z^{(n)}, \mathbf{x}^{(n)}; \theta)}{q_n(z^{(n)})} \right] \right) \end{aligned}$$

- We introduce auxilliary distributions $q_n(z^{(n)})$ over each of the latent variables

$$\begin{aligned}
 \sum_{n=1}^N \log \left(\sum_{z^{(n)}} p(z^{(n)}, \mathbf{x}^{(n)}; \boldsymbol{\theta}) \right) &= \sum_{n=1}^N \log \left(\sum_{z^{(n)}} q_n(z^{(n)}) \frac{p(z^{(n)}, \mathbf{x}^{(n)}; \boldsymbol{\theta})}{q_n(z^{(n)})} \right) \\
 &= \sum_{n=1}^N \log \left(\mathbb{E}_{q_n(z^{(n)})} \left[\frac{p(z^{(n)}, \mathbf{x}^{(n)}; \boldsymbol{\theta})}{q_n(z^{(n)})} \right] \right) \\
 &\geq \sum_{n=1}^N \mathbb{E}_{q_n(z^{(n)})} \left[\log \frac{p(z^{(n)}, \mathbf{x}^{(n)}; \boldsymbol{\theta})}{q_n(z^{(n)})} \right]
 \end{aligned}$$

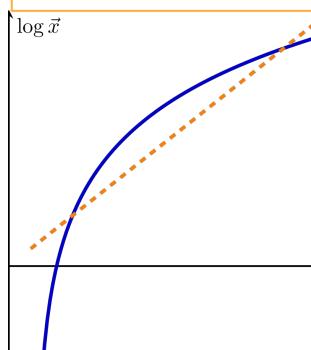
If g is **concave**, the inequality changes directions:

$$g(\mathbb{E}[X]) \leq \mathbb{E}[g(X)]$$

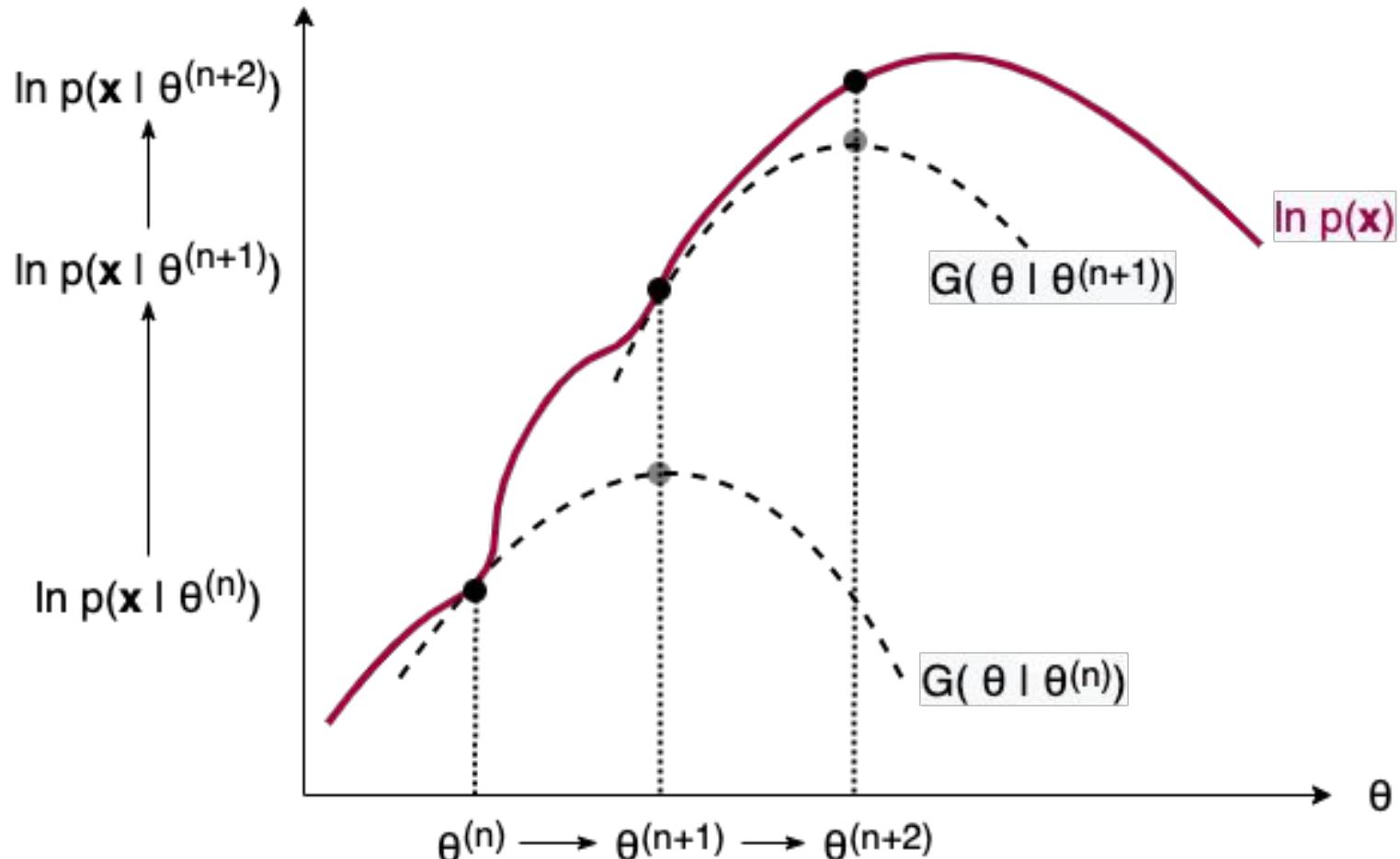
- In the last step, we use Jensen's Inequality. Since \log is concave:

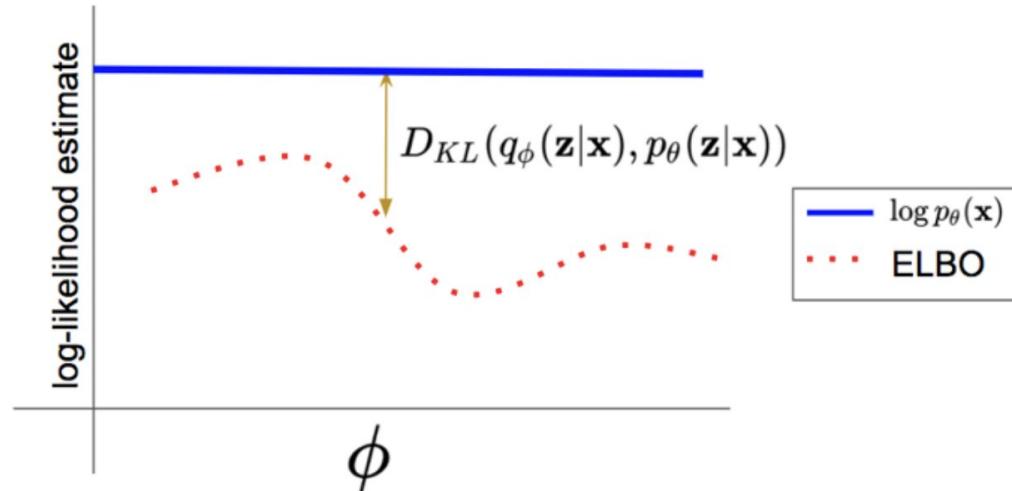
$$\log \left(\mathbb{E}_{q_n(z^{(n)})} \left[\frac{p(z^{(n)}, \mathbf{x}^{(n)}; \boldsymbol{\theta})}{q_n(z^{(n)})} \right] \right) \geq \mathbb{E}_{q_n(z^{(n)})} \left[\log \frac{p(z^{(n)}, \mathbf{x}^{(n)}; \boldsymbol{\theta})}{q_n(z^{(n)})} \right]$$

Called Evidence Lower Bound (**ELBO**)



Log-likelihood





$$\log p(\mathbf{x}; \theta) \geq \sum_{\mathbf{z}} q(\mathbf{z}; \phi) \log p(\mathbf{z}, \mathbf{x}; \theta) + H(q(\mathbf{z}; \phi)) = \underbrace{\mathcal{L}(\mathbf{x}; \theta, \phi)}_{\text{ELBO}}$$

$$\log p(\mathbf{x}; \theta) = \mathcal{L}(\mathbf{x}; \theta, \phi) + D_{KL}(q(\mathbf{z}; \phi) \| p(\mathbf{z}|\mathbf{x}; \theta))$$

The better $q(\mathbf{z}; \phi)$ can approximate the posterior $p(\mathbf{z}|\mathbf{x}; \theta)$, the smaller $D_{KL}(q(\mathbf{z}; \phi) \| p(\mathbf{z}|\mathbf{x}; \theta))$ we can achieve, the closer ELBO will be to $\log p(\mathbf{x}; \theta)$. Next: jointly optimize over θ and ϕ to maximize the ELBO over a dataset

- **Evidence lower bound (ELBO) holds for any q**

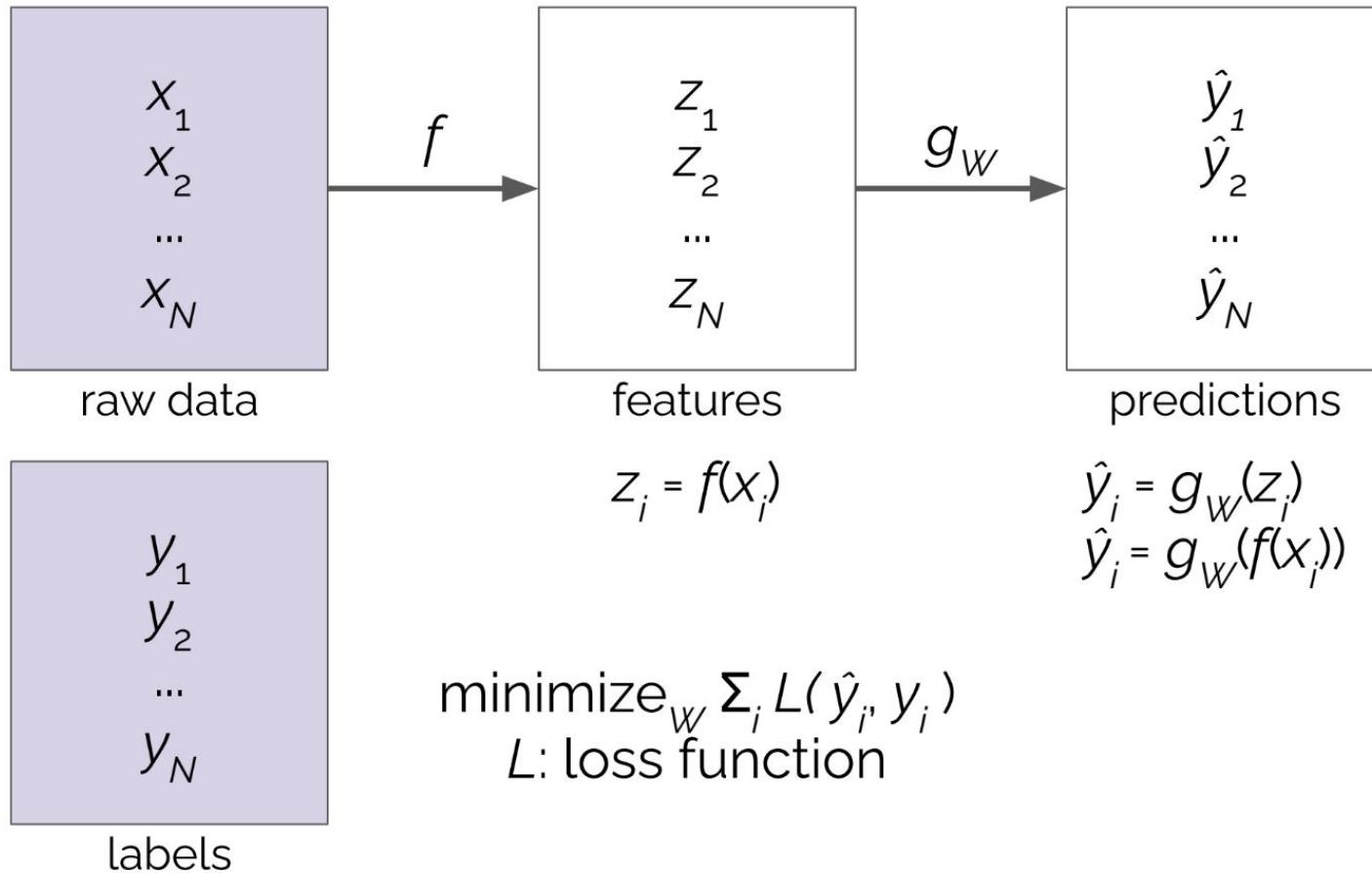
$$\log p(\mathbf{x}; \theta) \geq \sum_{\mathbf{z}} q(\mathbf{z}) \log p(\mathbf{z}, \mathbf{x}; \theta) + H(q)$$

- Equality holds if $q = p(\mathbf{z}|\mathbf{x}; \theta)$ because $D_{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}; \theta)) = 0$

$$\log p(\mathbf{x}; \theta) = \sum_{\mathbf{z}} q(\mathbf{z}) \log p(\mathbf{z}, \mathbf{x}; \theta) + H(q)$$

- Confirms our intuition that we seek likely completions \mathbf{z} given the observed values (evidence) \mathbf{x} .

Recap | Supervised Learning



Recap | Supervised Learning

"feature engineering"

When Does Contrastive Visual Representation Learning Work?

Elijah Cole¹ Xuan Yang² Kimberly Wilber² Oisin Mac Aodha^{2,4} Segei Belongie³
¹Cafeck ²Google ³University of Edinburgh ⁴Alan Turing Institute ⁵University of Copenhagen

Recent self-supervised representation learning techniques have largely closed the gap between supervised and unsupervised learning on ImageNet classification. While the particulars of pretraining on ImageNet are now relatively well understood, it is less clear what general best practices for replicating this success on other datasets. As a first step in this direction, we study contrastive self-supervised learning on four challenging visual domains. By looking through the lens of data quantity, data distribution, and task granularities, we provide insights into the necessary conditions for successful self-supervised learning. Our key findings include observations that (i) contrastive learning on ImageNet with just 500k images is modest, (ii) adding pretraining images from another domain does not lead to more general representations, (iii) contrastive learning is robust to different input on supervised and self-supervised pretraining, and (iv) contrastive learning tags for behind supervised learning on fine-grained visual classification tasks.

1. Introduction
Self-supervised learning (SSL) techniques can now produce visual representations which are competitive with representations generated by fully supervised networks for most computer vision tasks. This is a major breakthrough for computer vision, as removing the need for large amounts of labels at training time has the potential to scale up the field significantly. The success of SSL for computer vision is currently too difficult or costly to obtain. However, with some limited exceptions, the vast majority of current state-of-the-art SSL methods are trained on large standard datasets like ImageNet [1]. As a result, we do not have a good understanding of how well these methods work when trained on smaller datasets.
Under what conditions do self-supervised contrastive representation learning methods produce "good" visual representations? This question is important for computer vision researchers because it adds to our understanding of SSL and highlights opportunities for new methods. This is also an important question for domain experts with limited resources who might be interested in applying SSL to real-world problems. With these objectives in mind, we attempt to answer the following questions:
(i) What is the impact of data quantity? How many unlabeled images do we need for pretraining, and when is it better to use a small dataset for pretraining rather than the need for linear classifier training or end-to-end fine-tuning on a downstream task? In which regimes do self-supervised representations benefit from larger datasets?
(ii) What is the impact of the pretraining domain? How well do self-supervised representations trained on one domain perform when used to initialize representations by combining datasets? Do different pretraining domains lead to complementary representations?
(iii) What is the impact of the input quality? How robust are self-supervised methods to training time image corruption such as reduced resolution, compression artifacts, or noise? Does SSL benefit from training on images with poor downstream performance on uncurrupted images?
(iv) What is the impact of task granularity? Does SSL

arXiv:2105.05837v2 [cs.CV] 4 Apr 2022

1



raw document
data

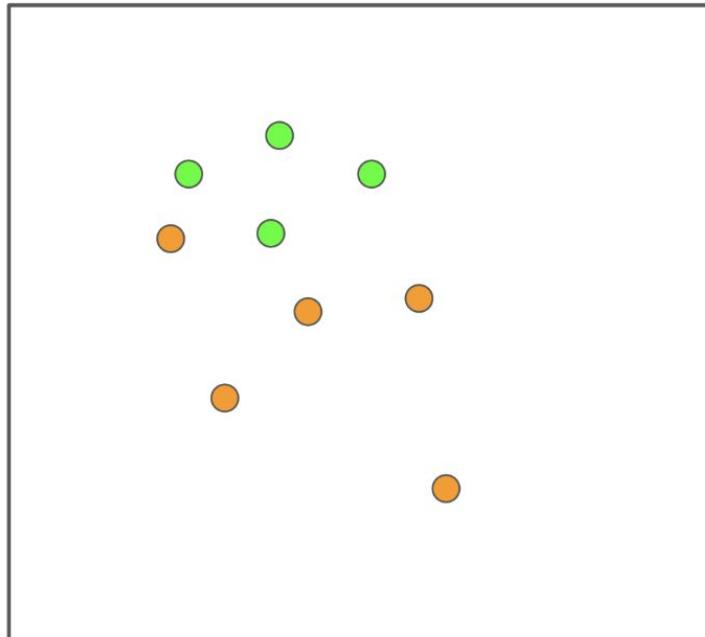
features

Word frequencies?

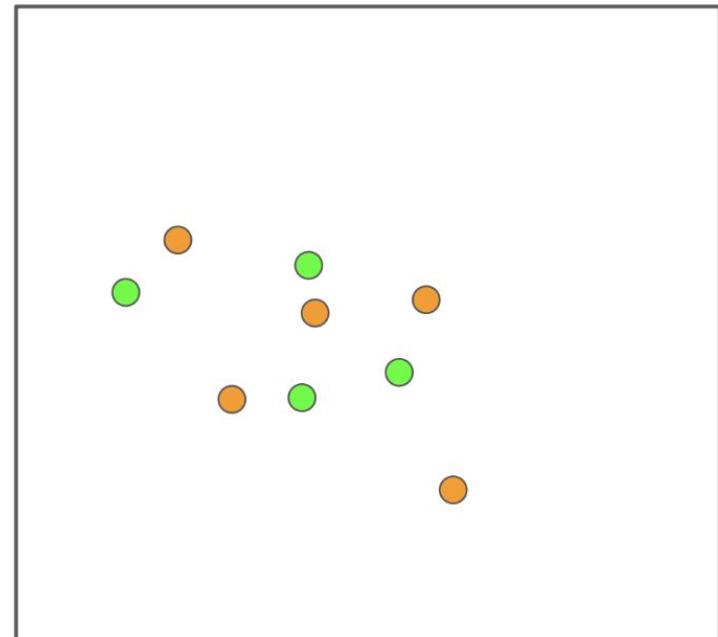
Word co-occurrences?

Section lengths?

good* features



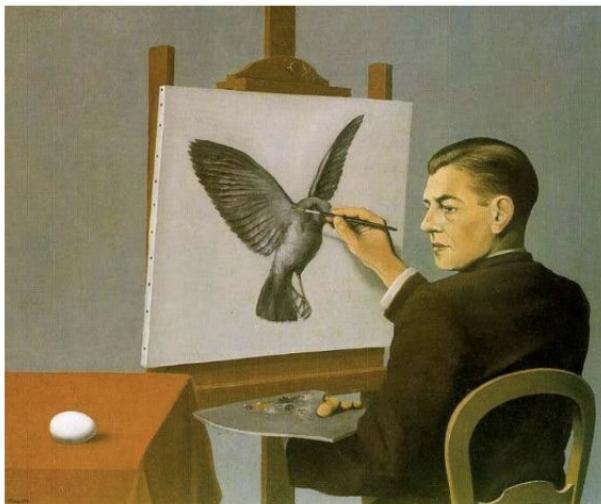
bad* features



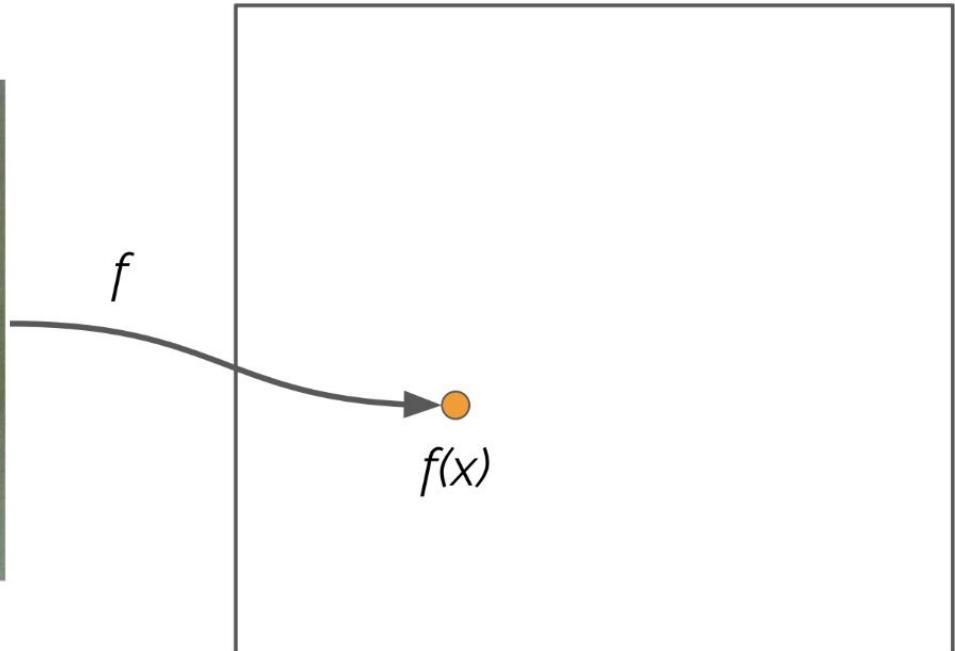
*For solving this binary classification task with a linear classifier.

Recap | Supervised Learning

It's hard to design good features, especially for complex objects.



X



Idea: Learn f from data.

In **deep learning**, we use special models (*deep neural networks*) which allow us to do *supervised learning* with *raw data* as inputs, eliminating the need for hand-designed features.

Example: Image Classification

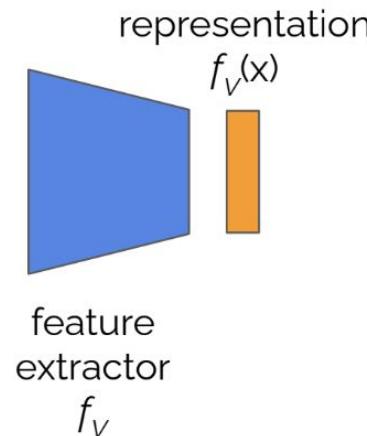


“Leopard
tortoise”

Example: Image Classification



image
 x



Example: Image Classification

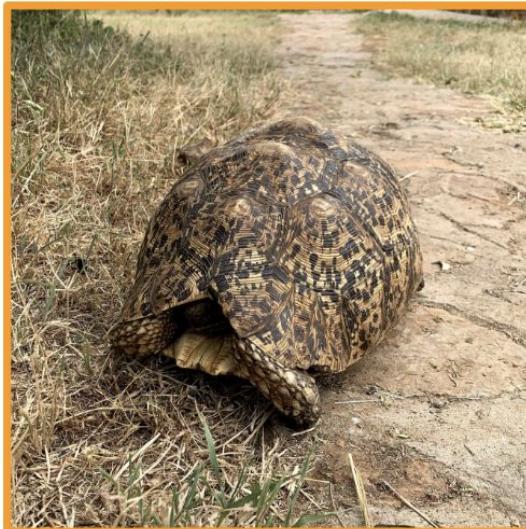
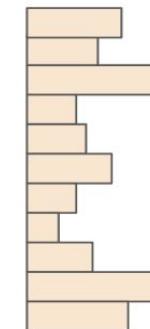
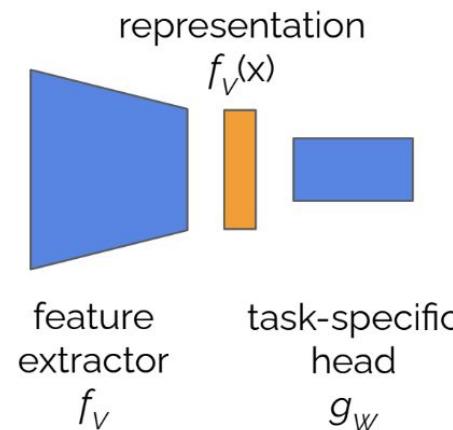


image
 x

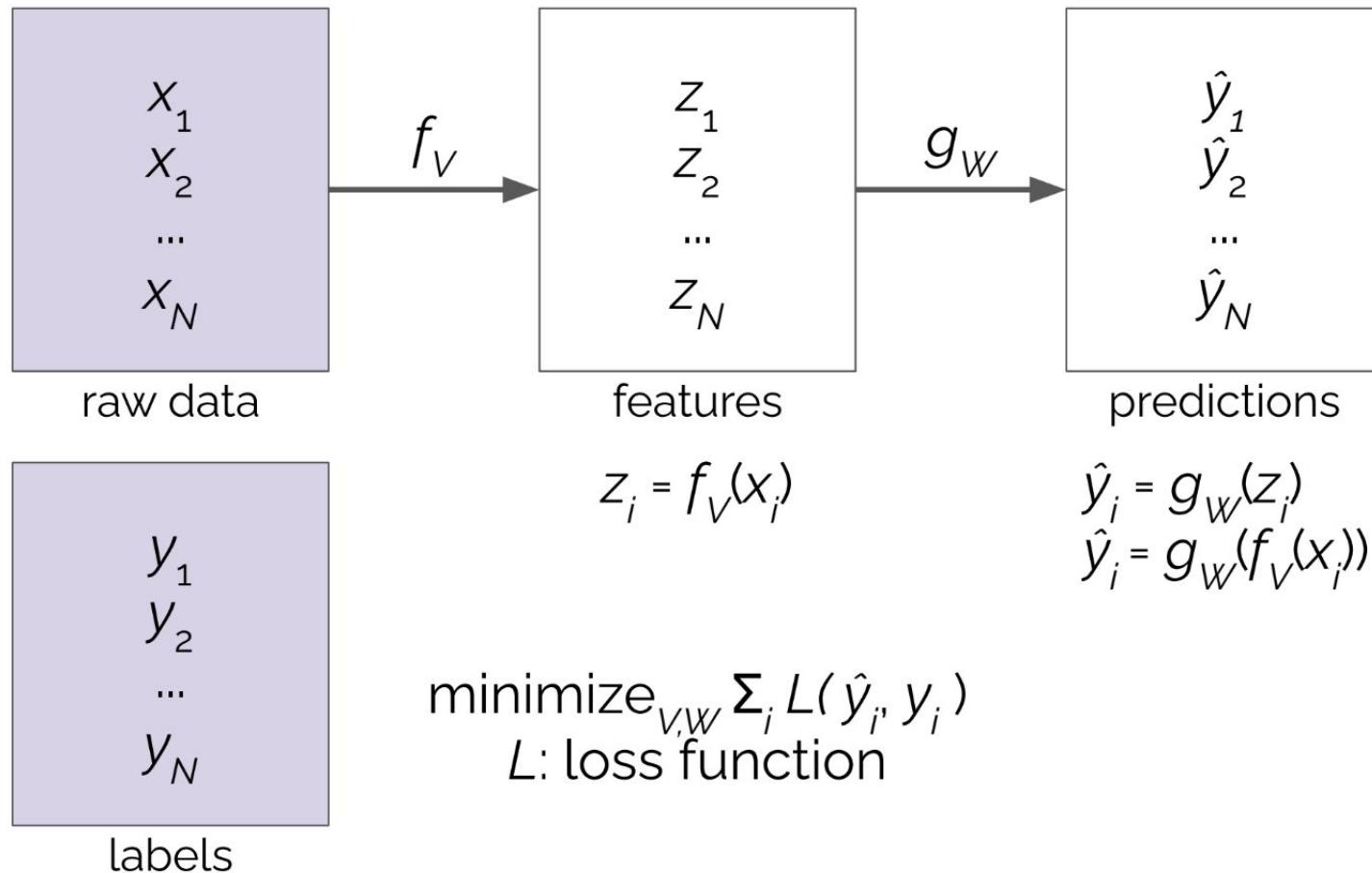


"African mud turtle"
"Painted dwarf gecko"
"Leopard tortoise"
"Helmeted chameleon"
...

$\hat{y}' = \text{softmax}(g_W(f_V(x)))$

$$V^*, W^* = \operatorname{argmin}_{V,W} E[L(\hat{y}, y)]$$

Recap | Supervised Learning



Recap | Supervised Learning

Summary:

- Supervised learning from hand-designed features
 - Feature engineering
- Supervised learning from raw data
 - AKA deep learning, representation learning
- Important ideas
 - Feature extractor + task-specific head
 - Task difficulty depends on features
 - Learning features from data is preferable
 - Lots of free parameters => need lots of data

Story so far:

- We know how to learn a representation from labeled data.
- We know how to re-use and evaluate a learned representation.
- **What if you don't have enough labeled data to learn the representation?**

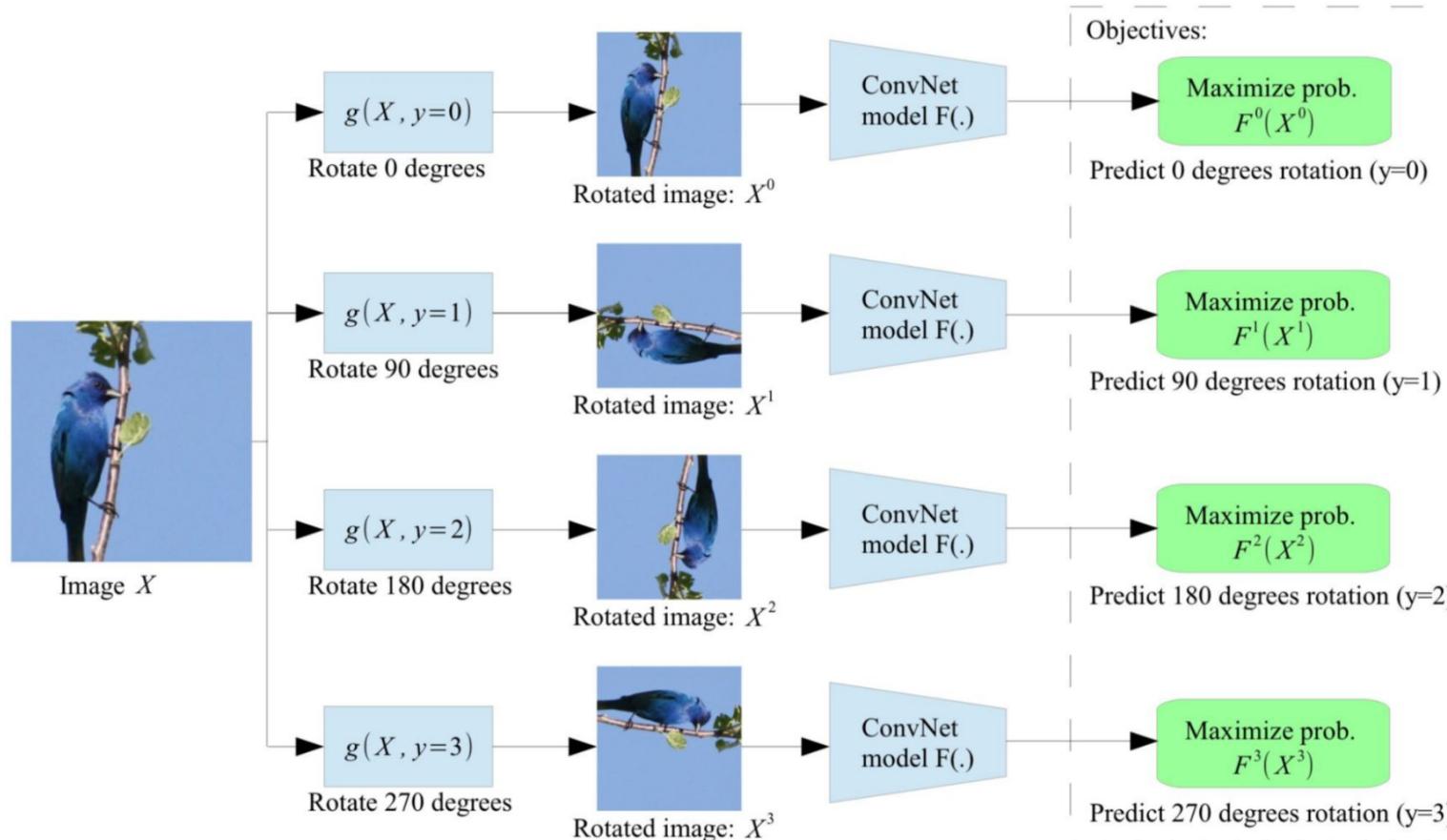
Self-Supervised (Unsupervised) Learning

Self-supervised pretraining looks a lot like supervised pretraining.

The only difference is that the labels are automatically generated *pseudo-labels* instead of human-curated ground truth labels. This is called a **pretext task**.

Then we do transfer learning.

SSL Basics | Pretext Task Examples



SSL Basics | Pretext Task Examples



(a) Input context



(b) Human artist



Inpainting ([Pathak et al. 2016](#))
(c) Context Encoder
(L_2 loss)

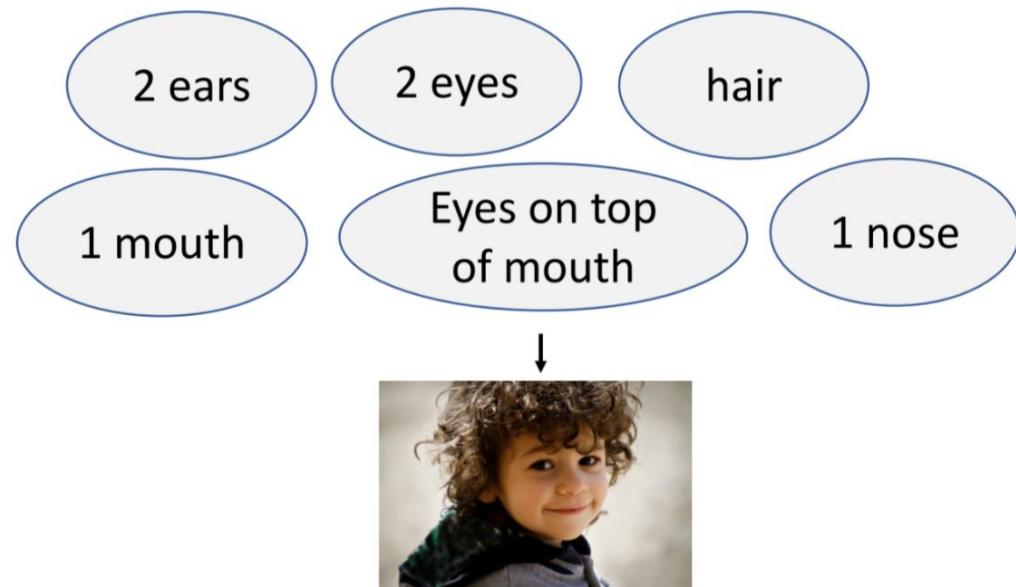
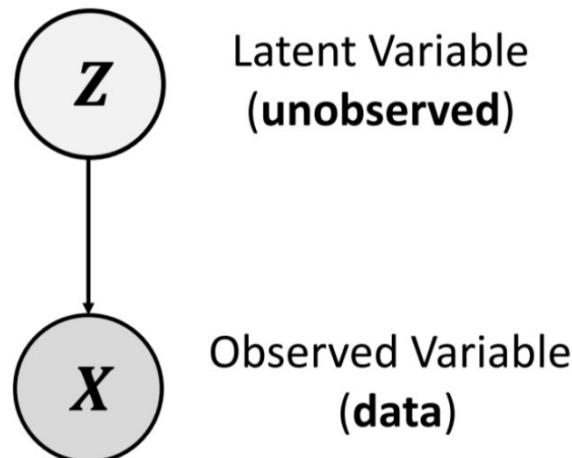


(d) Context Encoder
(L_2 + Adversarial loss)

(Variational) AutoEncoder

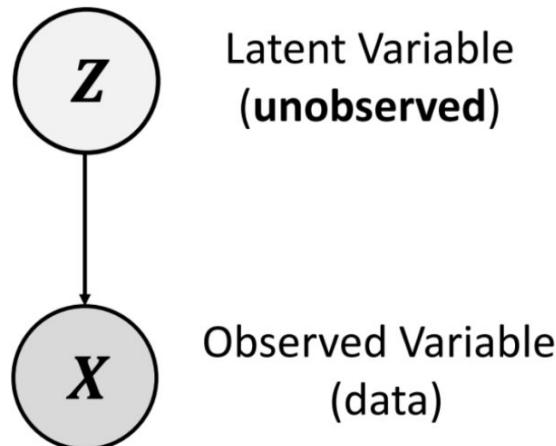
Variational Autoencoders

Belong in the family of latent variable probabilistic models that can infer the hidden structure in the underlying data



Variational Autoencoders

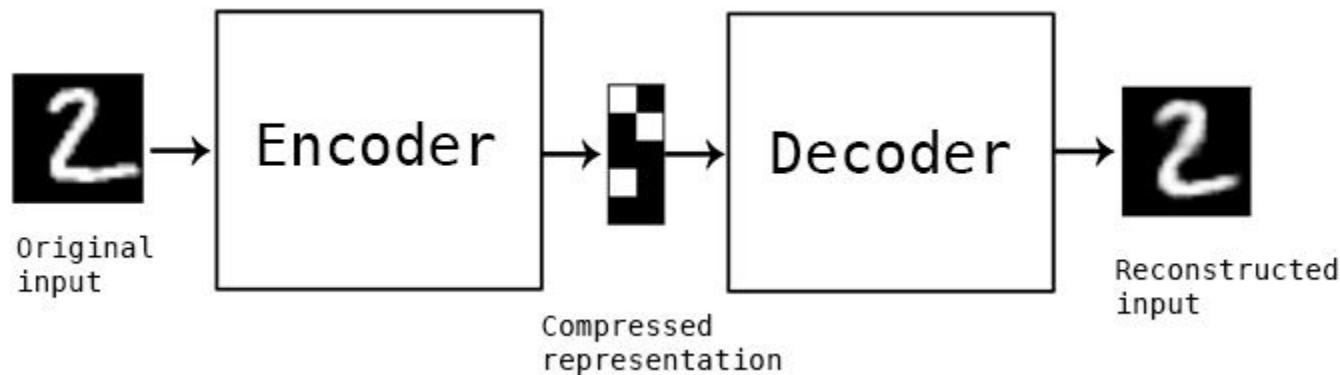
Belong in the family of latent variable probabilistic models that can infer the hidden structure in the underlying data



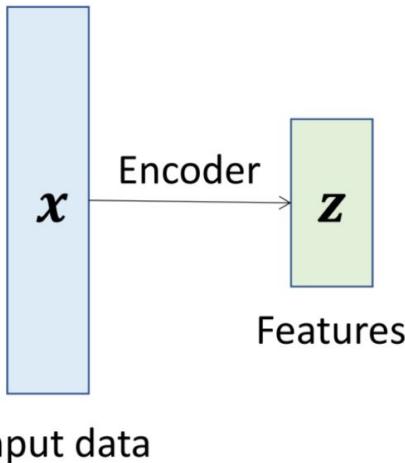
The goal of variational autoencoders is twofold:

1. Model the data distribution $p(x)$ – as is the case with all generative models
2. Extract a latent representation \mathbf{Z} that “describes” the data.

~~Variational Autoencoders~~

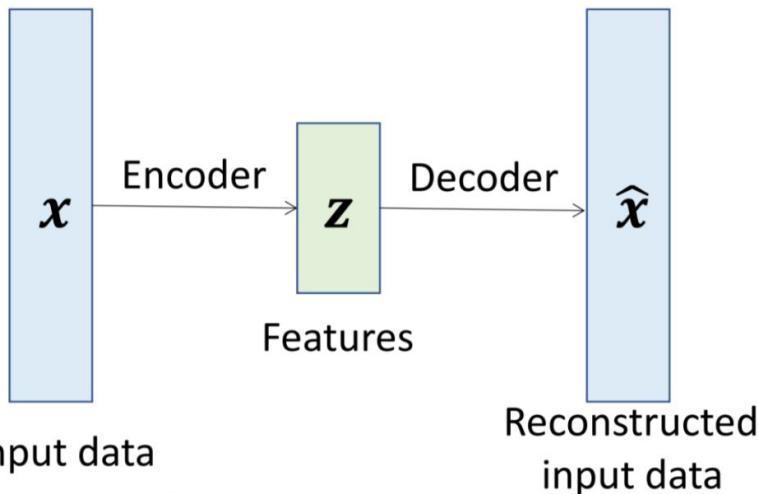


~~Variational Autoencoders~~



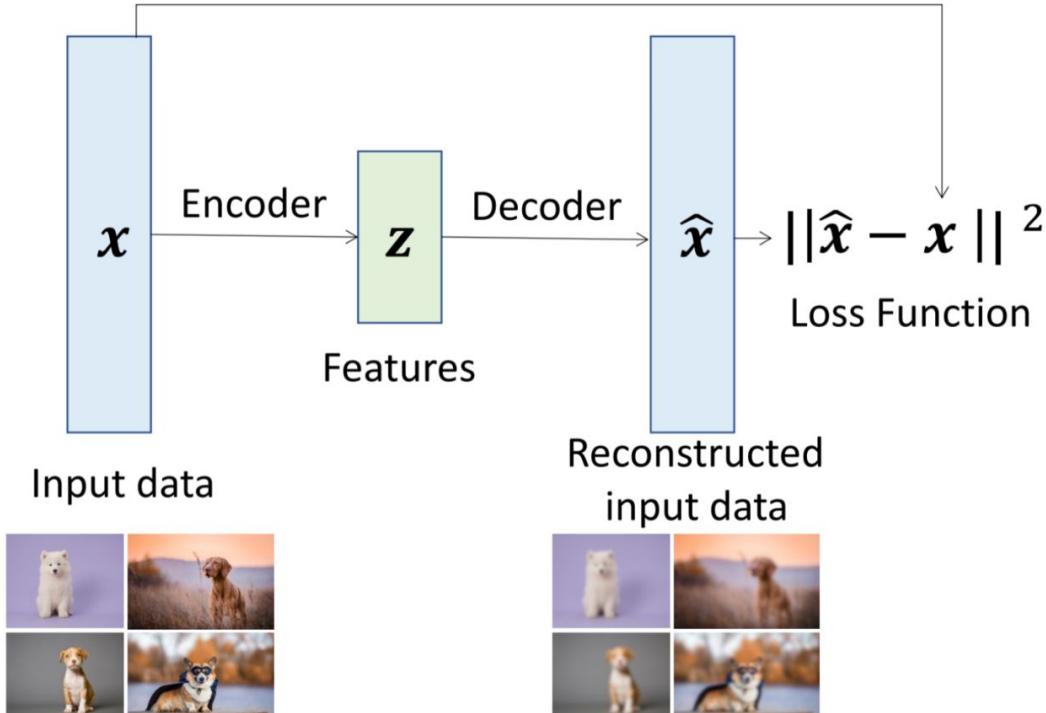
- The encoder extracts features which contain useful information from the data and can be used for downstream tasks
- We don't have labels, only data. So, we want to use unsupervised methods to learn these features.
- How?

~~Variational Autoencoders~~



- Use the features z to reconstruct the input data x with a decoder

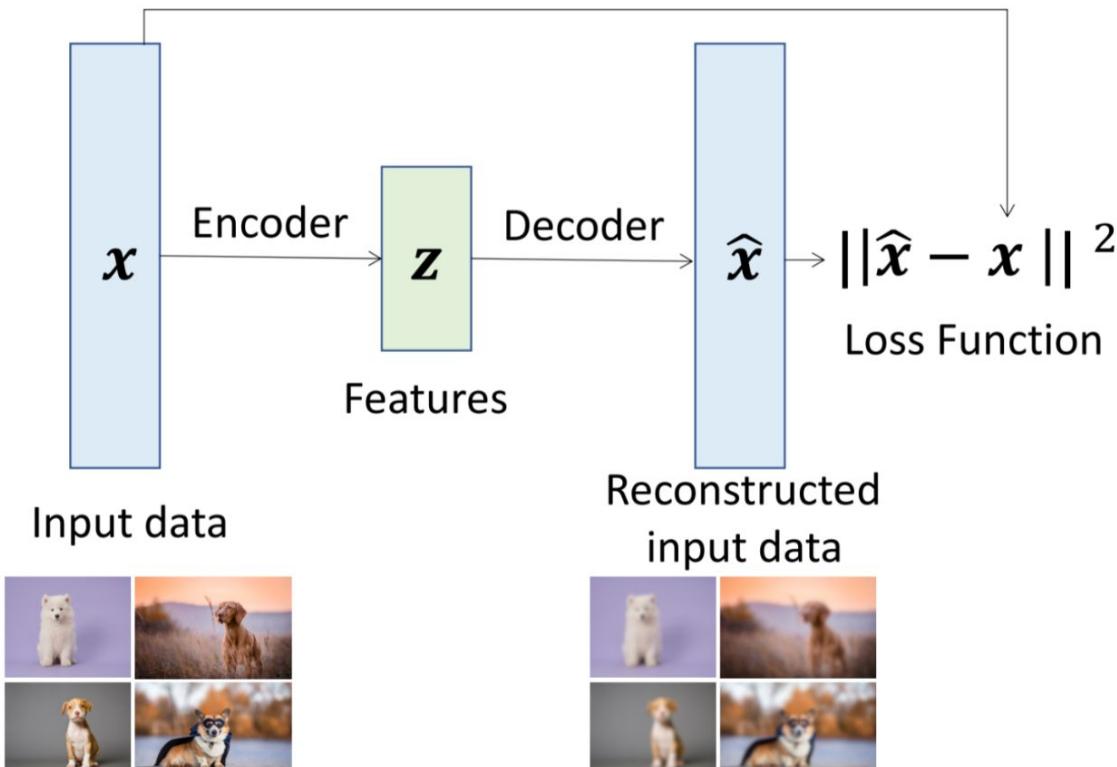
~~Variational Autoencoders~~



During training:

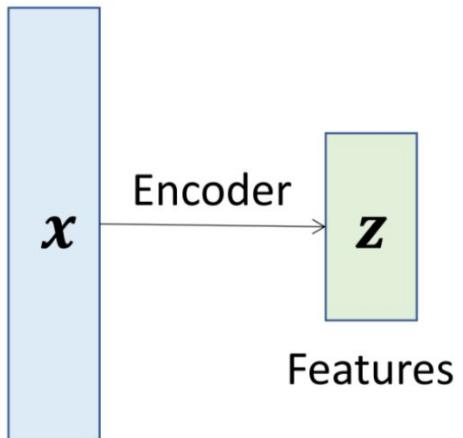
- We minimize an L2 distance between the input and the reconstructed data
- No labels used. Just data!

Autoencoders



The features z need to be lower dimensional than the data
(\rightarrow dimensionality reduction, data compression)

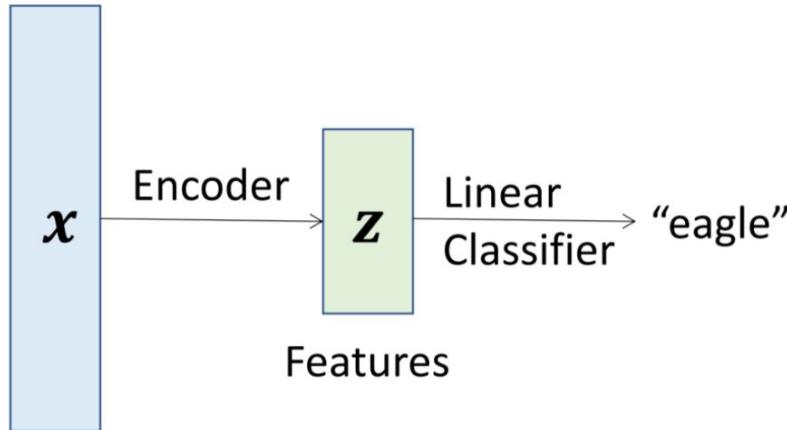
Autoencoders



After training:

- Throw away the decoder
- Use encoder to encode input images to “useful” features
- These features can be used for downstream tasks because they encode useful bits about the input (since one can reconstruct the input from them)

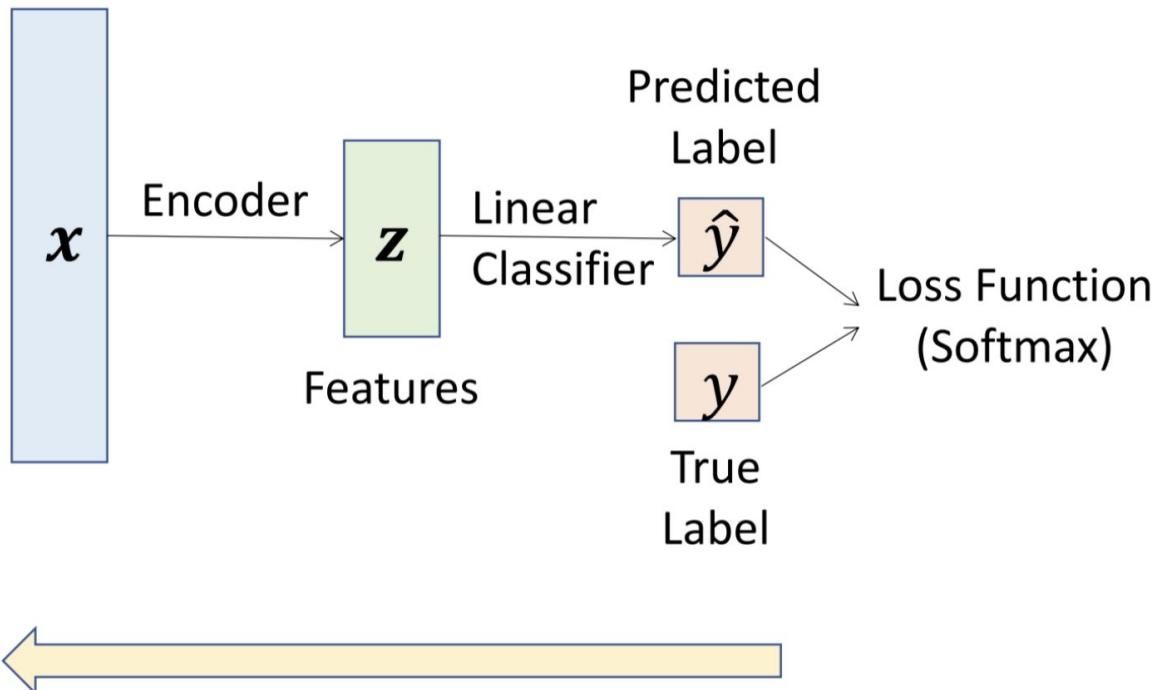
Autoencoders: Downstream Tasks



- A linear classifier inputs the features of the input image
- We train on a downstream supervised task with few annotated data points

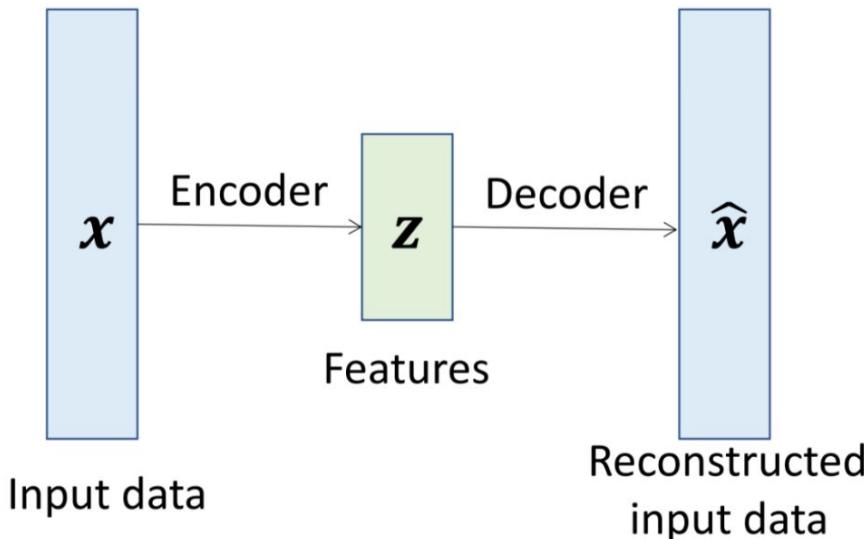


Autoencoders: Downstream Tasks



Finetune the encoder and the linear classifier

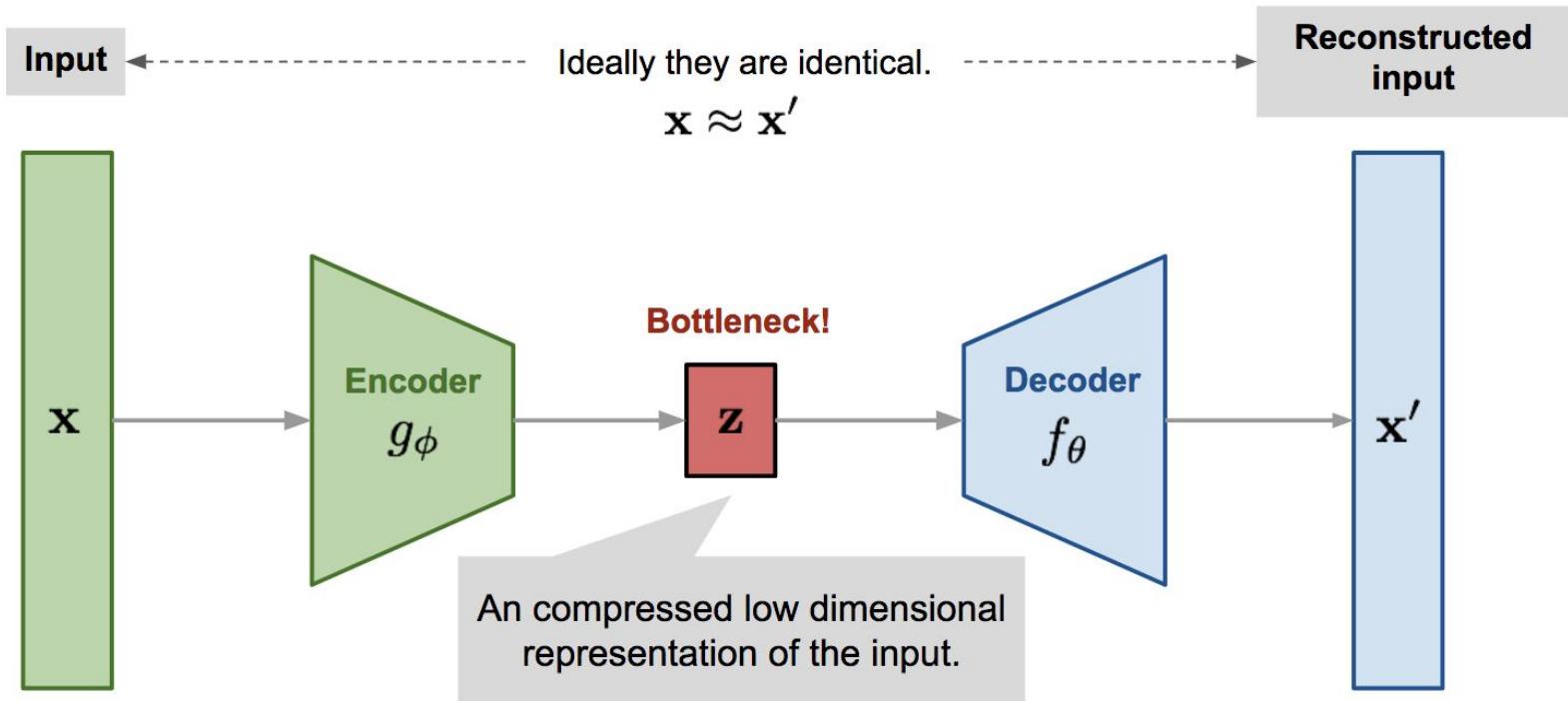
Autoencoders



✓ Autoencoders learn latent features for data without any labels

✓ The features can be used to initialize a supervised task

But they are not probabilistic and cannot be used to sample new data from. Remember we want to capture $p(x)$ and sample from it.

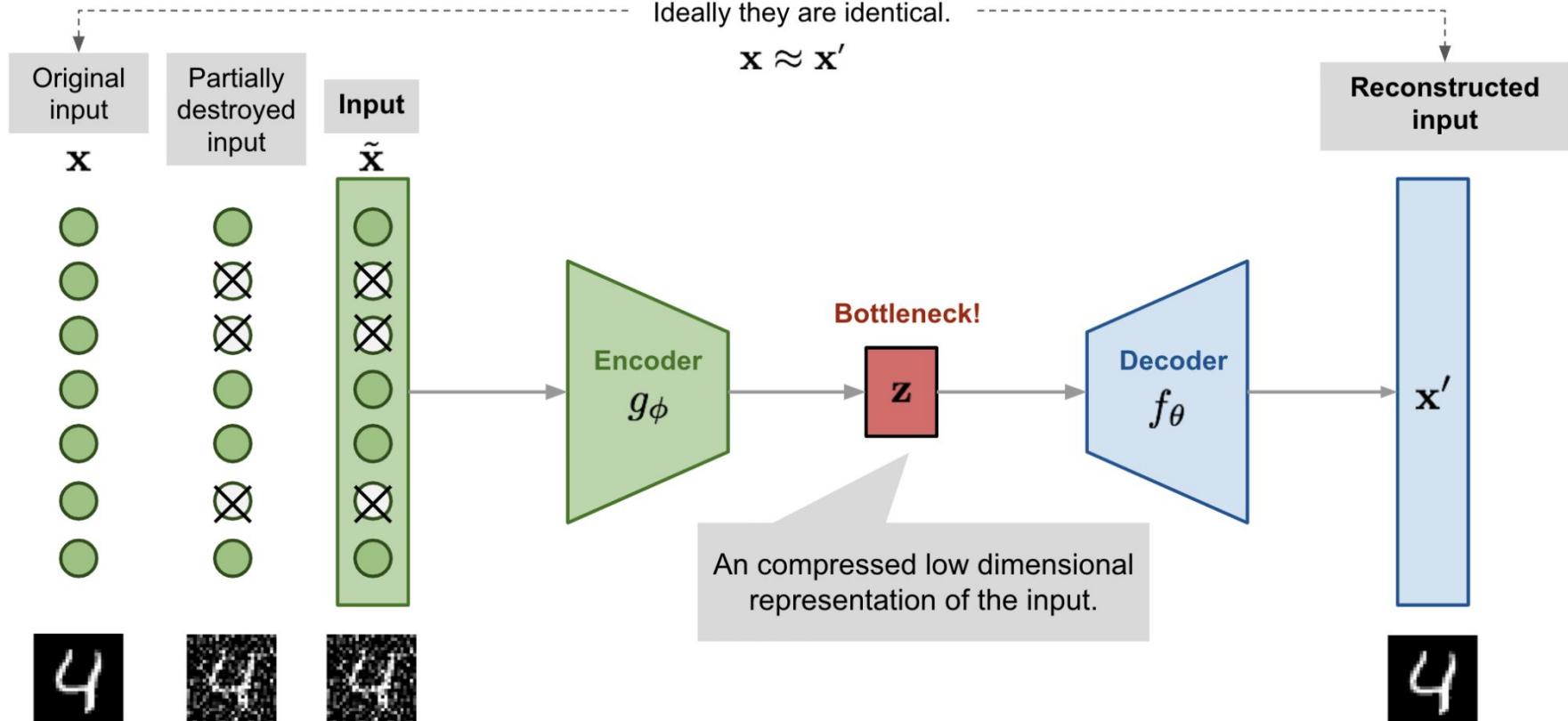


4

$$L_{\text{AE}}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - f_\theta(g_\phi(\mathbf{x}^{(i)})))^2$$

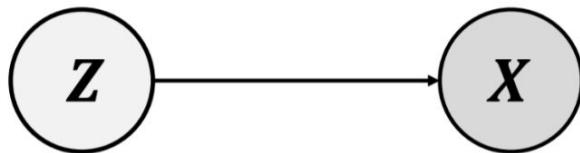
4

Denoising Auto-Encoder



Variational Autoencoders

Assume we have training data $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$ which are generated by an unobserved (latent) representation \mathbf{z} .

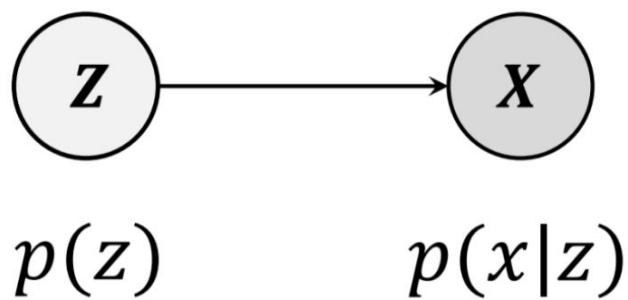


For example, \mathbf{x} is a dog image, and \mathbf{z} is a latent vector which captures all the information one might need to generate an image of a dog, e.g. it has 2 eyes, 2 ears, 4 legs, it's cute etc.

Variational Autoencoders

The joint distribution

$$p(x, z) = p(x|z)p(z)$$



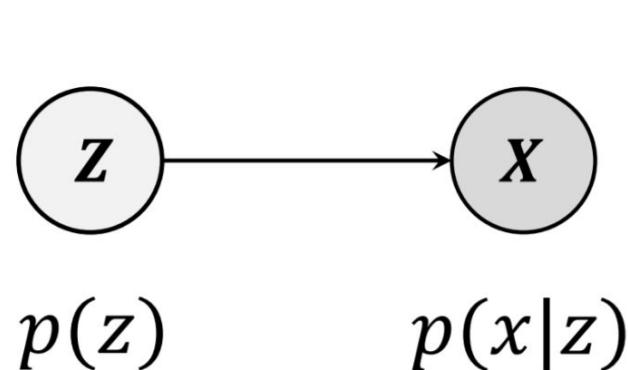
After training,

1. Sample from prior: $\hat{z} \sim p(z)$
2. Sample data from conditional: $\hat{x} \sim p(x|\hat{z})$

Intuitively, first generate the “high-level” semantic information about the data. Then generate the data.

Variational Autoencoders

Assume a simple prior $p(z)$, e.g. Gaussian



Represent $p(x|z)$ with a neural network

- The conditional distribution is parametrized, $p_\theta(x|z)$
- This neural network is the decoder

Starting from Gaussian, we can approximate very complicated data distribution (even real-world images / videos)

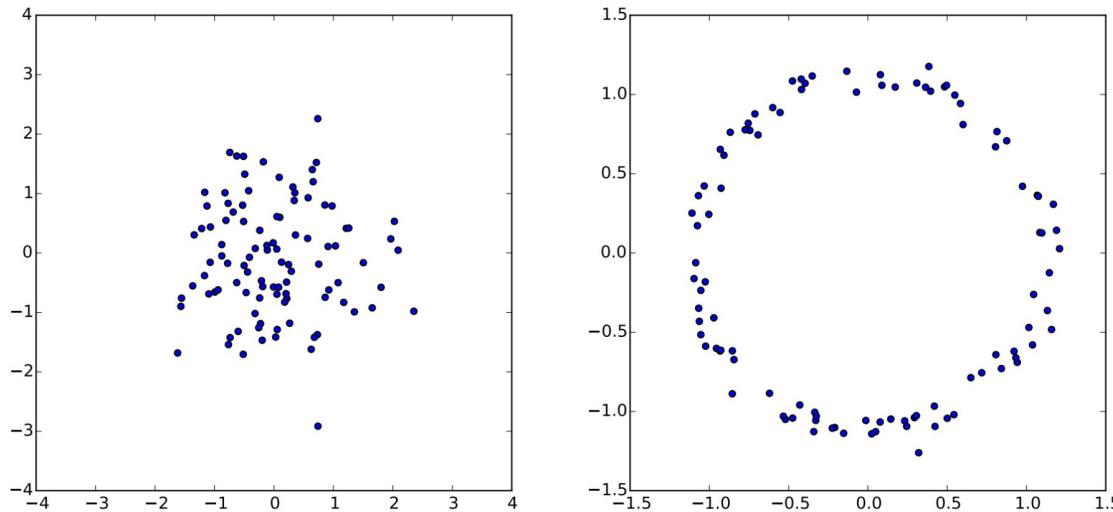
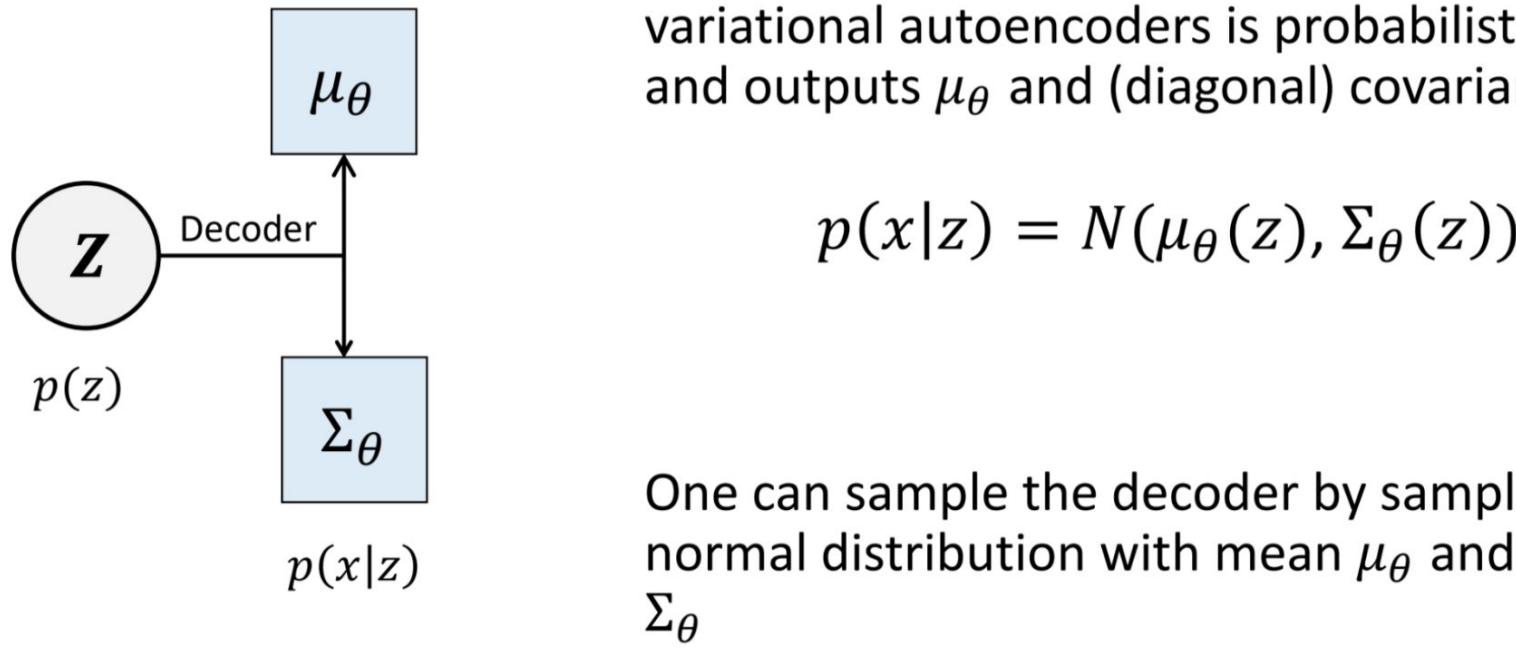


Figure 2: Given a random variable z with one distribution, we can create another random variable $X = g(z)$ with a completely different distribution. Left: samples from a gaussian distribution. Right: those same samples mapped through the function $g(z) = z/10 + z/\|z\|$ to form a ring. This is

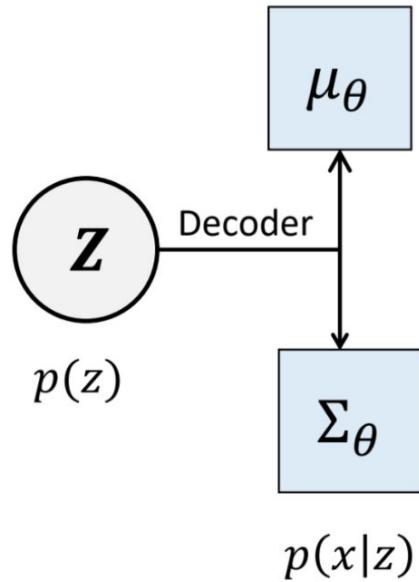
Variational Autoencoders

Unlike in original autoencoders, the decoder in variational autoencoders is probabilistic. It inputs z and outputs μ_θ and (diagonal) covariance Σ_θ



Variational Autoencoders

Density Estimation: We want to be able to estimate the density $p_\theta(x)$ of each datapoint x .

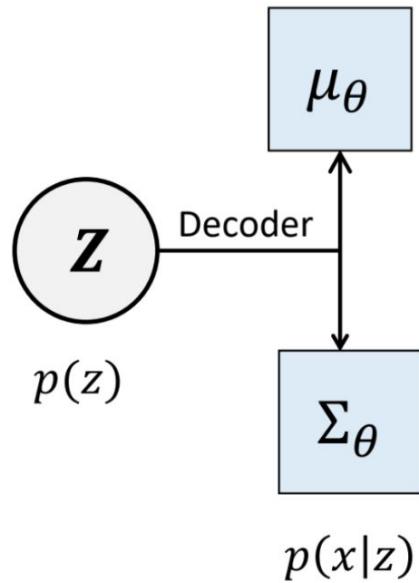


Note that we don't observe z , so to get the density we need to marginalize

$$p_\theta(x) = \int p_\theta(x, z) dz = \int p_\theta(x|z) p(z) dz$$

Variational Autoencoders

Density Estimation: We want to be able to estimate the density $p_\theta(x)$ of each datapoint x .



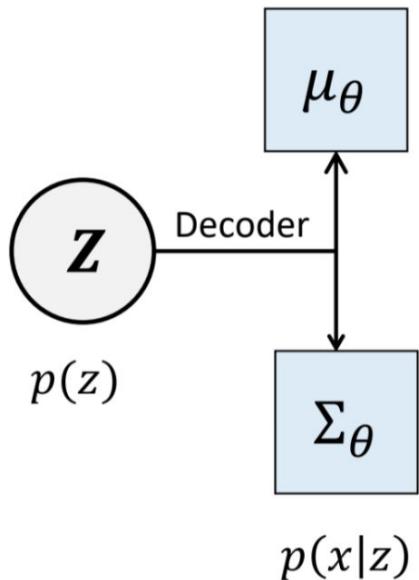
Note that we don't observe \mathbf{z} , so to get the density we need to marginalize

$$p_\theta(x) = \int p_\theta(x, z) dz = \int p_\theta(x|z) p(z) dz$$

Impossible to integrate
over all z . Intractable!

Variational Autoencoders

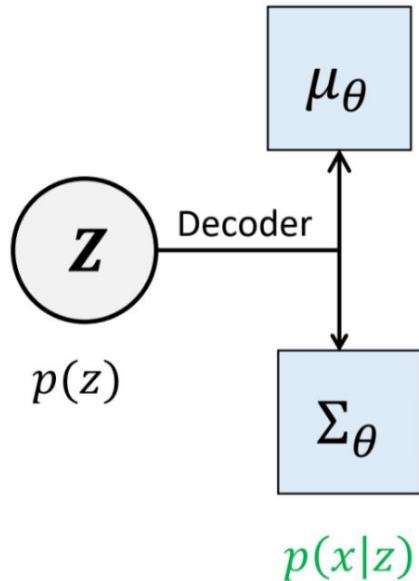
Density Estimation: We want to be able to estimate the density $p_\theta(x)$ of each datapoint x .



Bayes' rule says:

$$p(x) = \frac{p(x|z)p(z)}{p(z|x)}$$

Variational Autoencoders



Density Estimation: We want to be able to estimate the density $p_\theta(x)$ of each datapoint x .

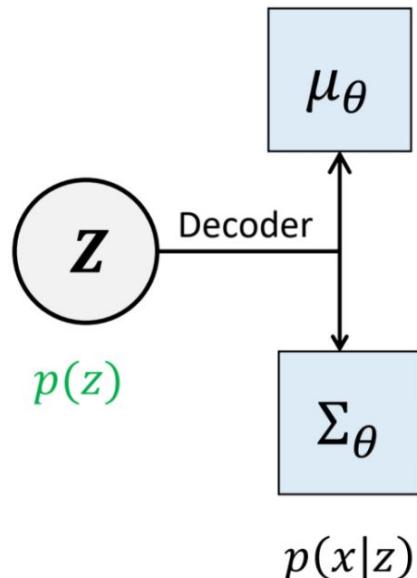
Bayes' rule says:

$$p(x) = \frac{p(x|z)p(z)}{p(z|x)}$$

Easy! My decoder network outputs $p(x|z)$ parametrized as a gaussian with $\mu_\theta, \Sigma_\theta$

Variational Autoencoders

Density Estimation: We want to be able to estimate the density $p_\theta(x)$ of each datapoint x .



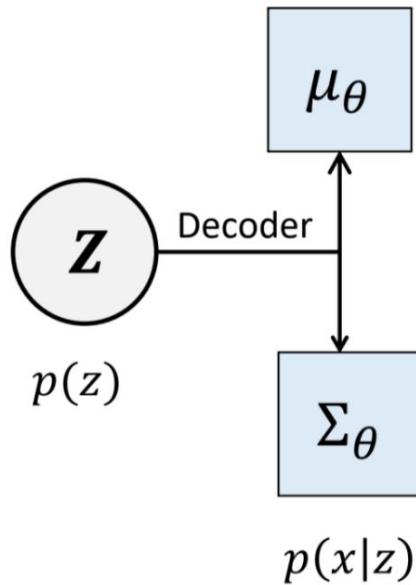
Bayes' rule says:

$$p(x) = \frac{p(x|z)p(z)}{p(z|x)}$$

Easy! I have assumed a Gaussian prior $N(0, I)$.

Variational Autoencoders

Density Estimation: We want to be able to estimate the density $p_\theta(x)$ of each datapoint x .



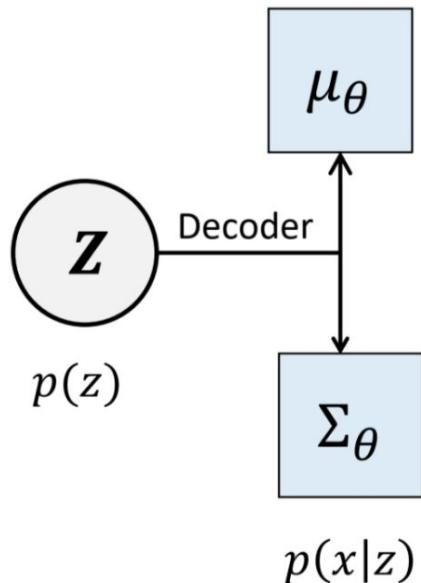
Bayes' rule says:

$$p(x) = \frac{p(x|z)p(z)}{p(z|x)}$$

Not easy! I have no way of computing this.

Variational Autoencoders

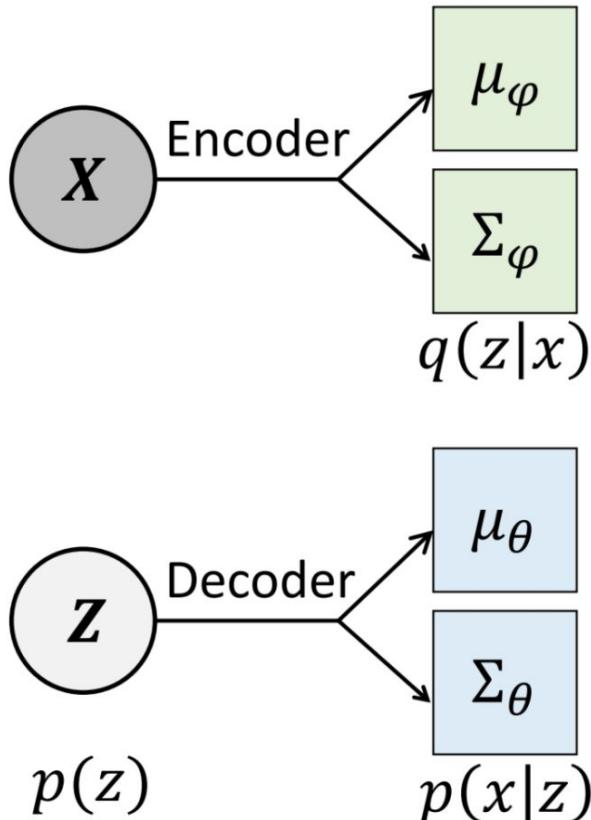
$$p(x) = \frac{p(x|z)p(z)}{p(z|x)}$$



We introduce a new parametric distribution $q(z|x)$ which will approximate the unknown $p(z|x)$

$$q(z|x) \approx p(z|x)$$

Variational Autoencoders



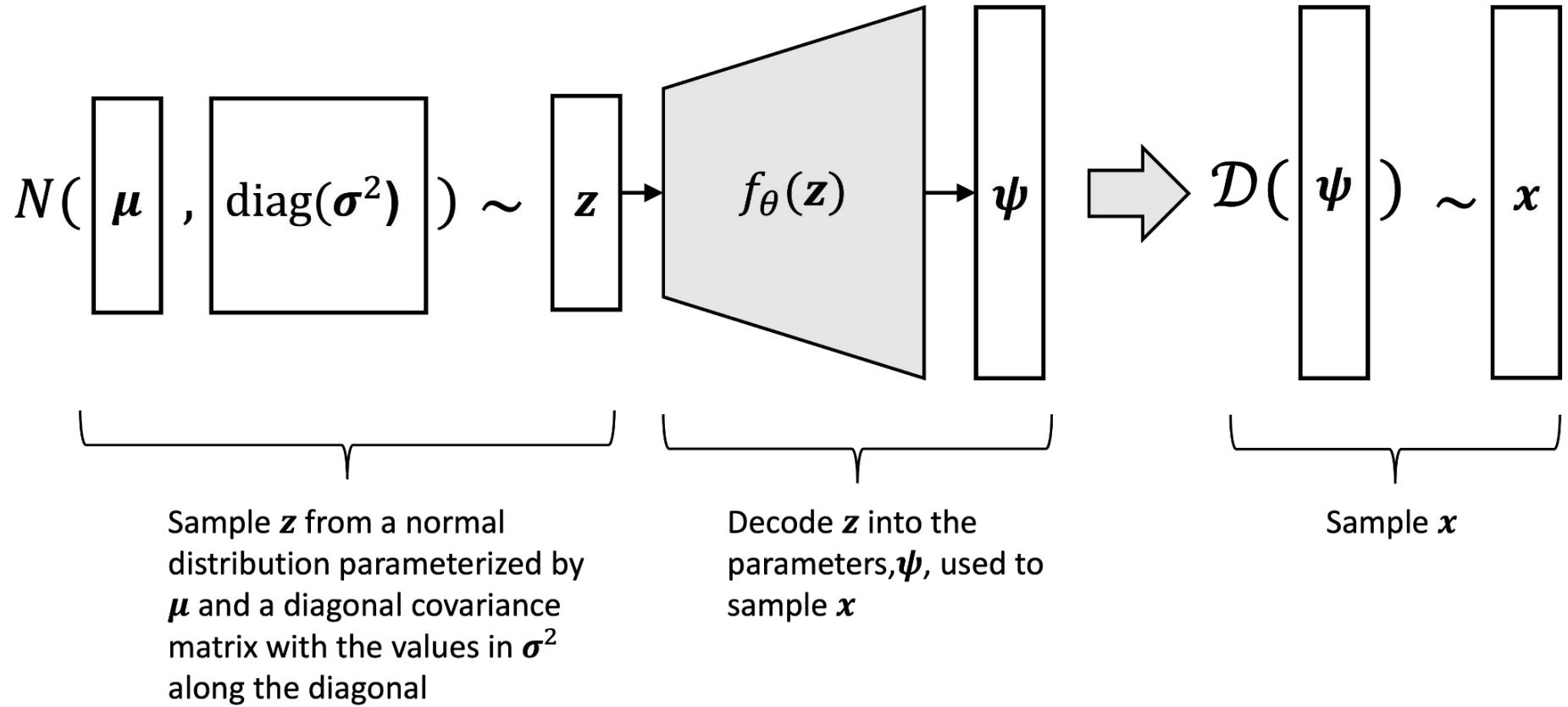
$$p(x) = \frac{p(x|z) p(z)}{p(z|x)} \approx \frac{p(x|z)p(z)}{q(z|x)}$$

$q(z|x)$ is parametrized by an encoder network with parameters φ which maps the input data to mean μ_φ and a **diagonal** covariance Σ_φ such that

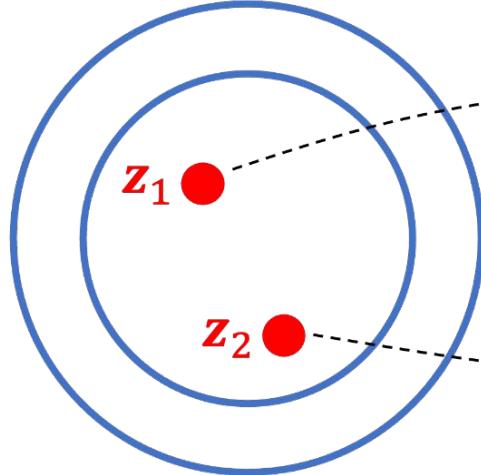
$$q_\varphi(z|x) = N(\mu_\varphi, \Sigma_\varphi)$$

Why diagonal Σ_φ ?

It makes computations easy.



$$p(\mathbf{z})$$



$$f_{\theta}(z_1)$$

$$f_{\theta}(z_2)$$

$$p(x | z_1)$$

$$x_1$$

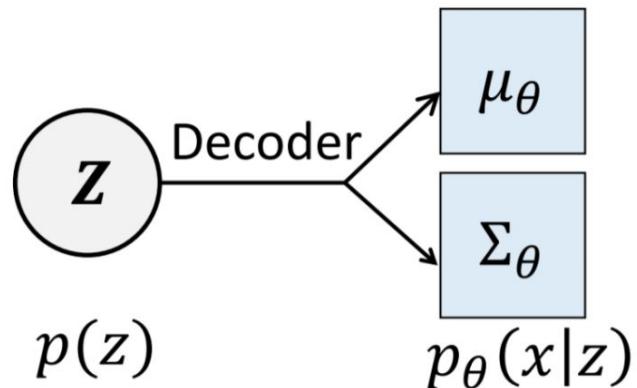
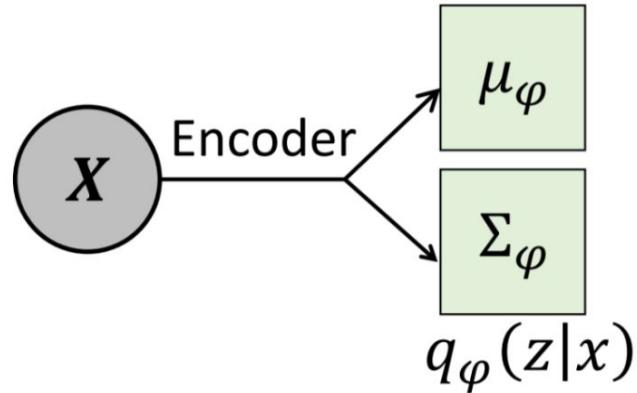
$$p(x | z_2)$$

$$x_2$$

$$\mathbb{R}^J$$

$$\mathbb{R}^D$$

Variational Autoencoders



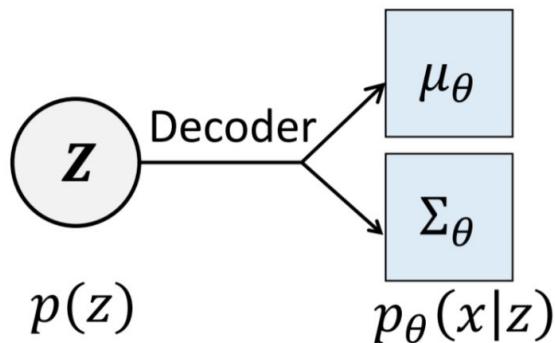
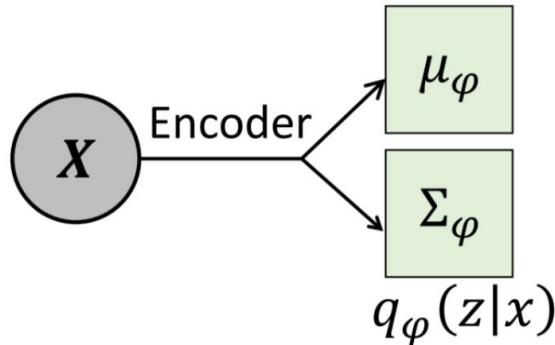
Optimization in VAEs involves jointly optimizing

1. The generative model parameters θ (decoder) to reduce the reconstruction error between input and output
2. The encoder parameters φ to model the posterior $q_\varphi(z|x)$

$$\max_{\theta, \varphi} \mathcal{L}_{\theta, \varphi}(x) = \mathbb{E}_{z \sim q_\varphi} [\log p_\theta(x|z)] - d_{KL}(q_\varphi(z|x) \| p(z))$$

Reconstruction Loss:

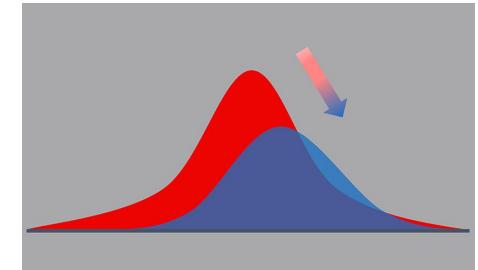
This term measures given encoded latent code (z), how well we can reconstruct data (x)



Encoder Prior Loss:

This term measures the distance between encoded latent distribution $q(z|x)$ to prior gaussian $p(z)$.

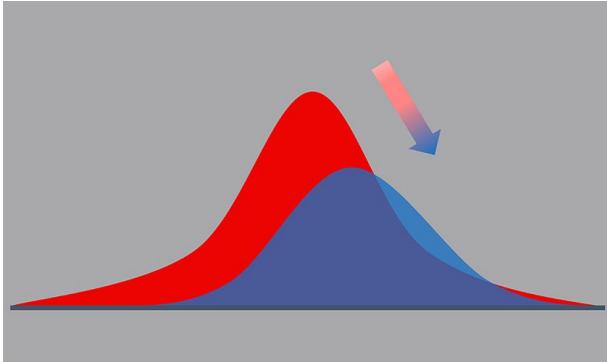
After this term becomes small



How to derive it?

Note that we want to learn $q(z|x)$ to approximate $p(z|x)$.

Let's start with the KL Divergence between them.



$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

$$D_{\text{KL}}(P \parallel Q) = \int_{-\infty}^{\infty} p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$

$$d_{KL}(q_\varphi(z|x) \mid\mid p_\theta(z|x)) = \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(z|x)]$$

Definition of KL divergence

$$\begin{aligned} d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(z|x)] \\ &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z) + \log p_\theta(x)] \end{aligned}$$

From Bayes' rule

$$\begin{aligned} d_{KL}(q_\varphi(z|x) \mid\mid p_\theta(z|x)) &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(z|x)] \\ &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z) + \log p_\theta(x)] \\ &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z)] + \log p_\theta(x) \end{aligned}$$

Since $p_\theta(x)$ does not depend on z we can move it out of the expectation

$$\begin{aligned} d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(z|x)] \\ &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z) + \log p_\theta(x)] \\ &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z)] + \textcolor{orange}{\log p_\theta(x)} \end{aligned}$$

Ah! This is exactly the objective we want to maximize! Maximum likelihood estimation (MLE)

$$\begin{aligned}
d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(z|x)] \\
&= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z) + \log p_\theta(x)] \\
&= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z)] + \log p_\theta(x)
\end{aligned}$$

\Rightarrow

$$\log p_\theta(x) - d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) = \mathbb{E}_{z \sim q_\varphi} [\log p_\theta(x|z)] - d_{KL}(q_\varphi(z|x) \mid p(z))$$

The definition of KL divergence

$$\begin{aligned} d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(z|x)] \\ &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z) + \log p_\theta(x)] \\ &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z)] + \log p_\theta(x) \end{aligned}$$

\Rightarrow

$$\log p_\theta(x) - d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) = \mathbb{E}_{z \sim q_\varphi} [\log p_\theta(x|z)] - d_{KL}(q_\varphi(z|x) \mid p(z))$$



The quantity we want to maximize!

$$\begin{aligned} d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(z|x)] \\ &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z) + \log p_\theta(x)] \\ &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z)] + \log p_\theta(x) \end{aligned}$$

\Rightarrow

$$\log p_\theta(x) - d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) = \mathbb{E}_{z \sim q_\varphi} [\log p_\theta(x|z)] - d_{KL}(q_\varphi(z|x) \mid p(z))$$


≥ 0 . We can't compute this but
this term will become small if the
encoder is high capacity. Why?

$$\begin{aligned}
d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(z|x)] \\
&= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z) + \log p_\theta(x)] \\
&= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z)] + \log p_\theta(x)
\end{aligned}$$

\Rightarrow

$$\log p_\theta(x) - d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) = \mathbb{E}_{z \sim q_\varphi} [\log p_\theta(x|z)] - d_{KL}(q_\varphi(z|x) \mid p(z))$$



This is a lower bound on our
desired objective: $\log p_\theta(x)$
The VAE objective!

$$\begin{aligned} d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(z|x)] \\ &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z) + \log p_\theta(x)] \\ &= \mathbb{E}_{z \sim q_\varphi} [\log q_\varphi(z|x) - \log p_\theta(x|z) - \log p(z)] + \log p_\theta(x) \end{aligned}$$

⇒

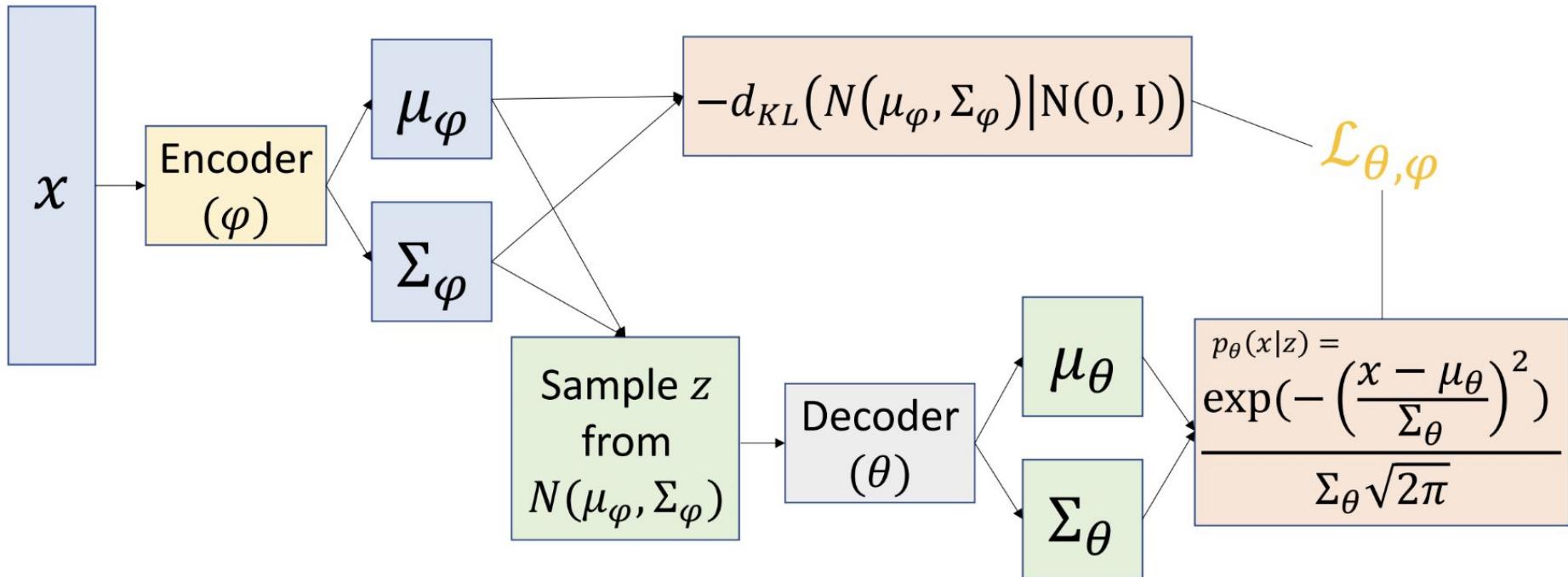
$$\log p_\theta(x) - d_{KL}(q_\varphi(z|x) \mid p_\theta(z|x)) = \mathbb{E}_{z \sim q_\varphi} [\log p_\theta(x|z)] - d_{KL}(q_\varphi(z|x) \mid p(z))$$



It's called the Evidence Lower
Bound (ELBO)

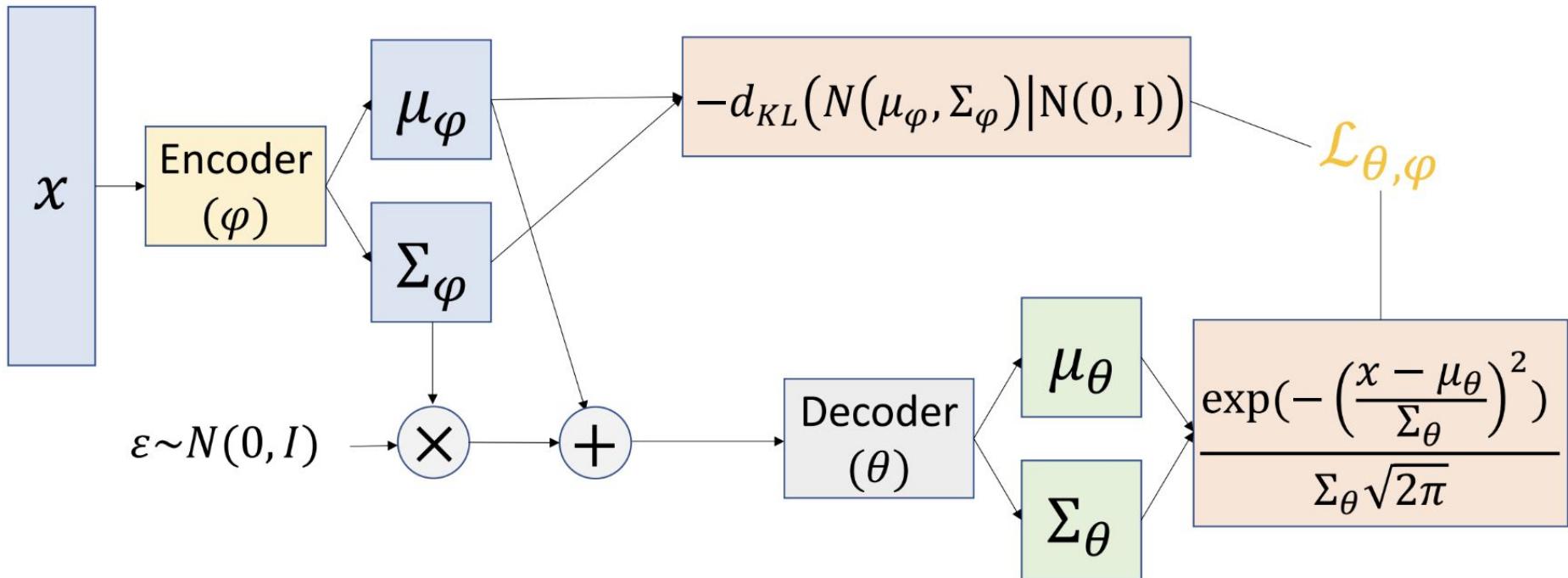
$$\mathcal{L}_{\theta,\varphi}(x) = \mathbb{E}_{z \sim q_{\varphi}} [\log p_{\theta}(x|z)] - d_{KL}(q_{\varphi}(z|x) \| p(z))$$

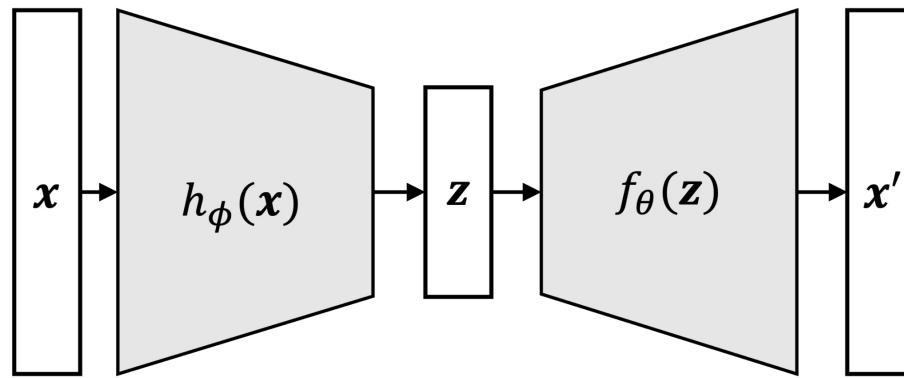
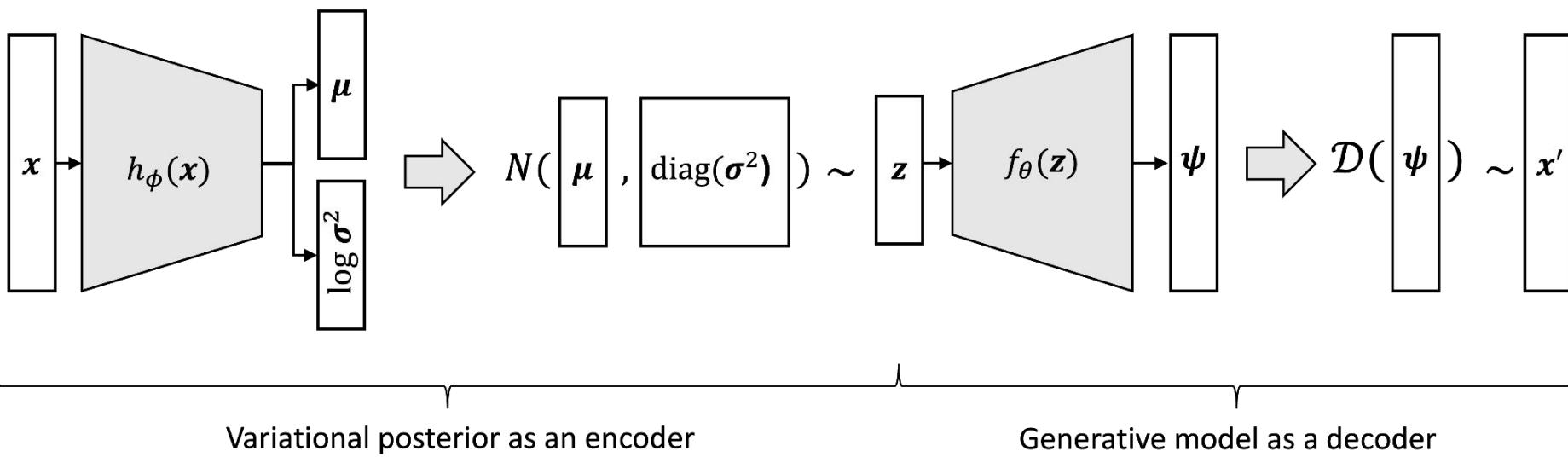
Training VAEs



Reparameterization trick

Training VAEs



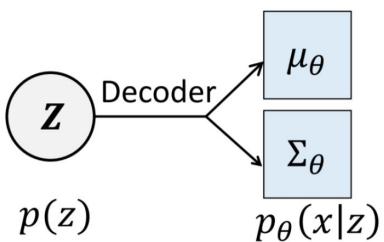
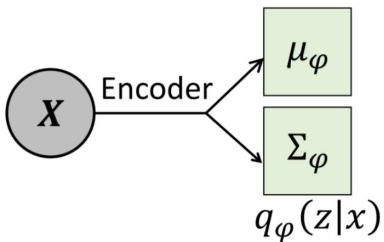


Optimizing the VAE objective

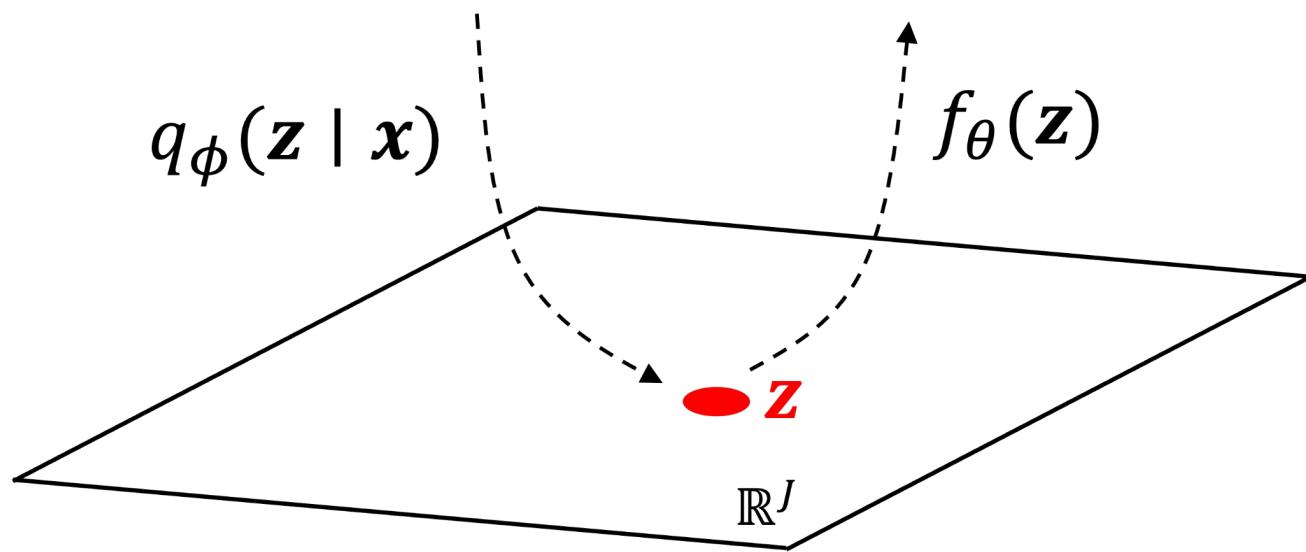
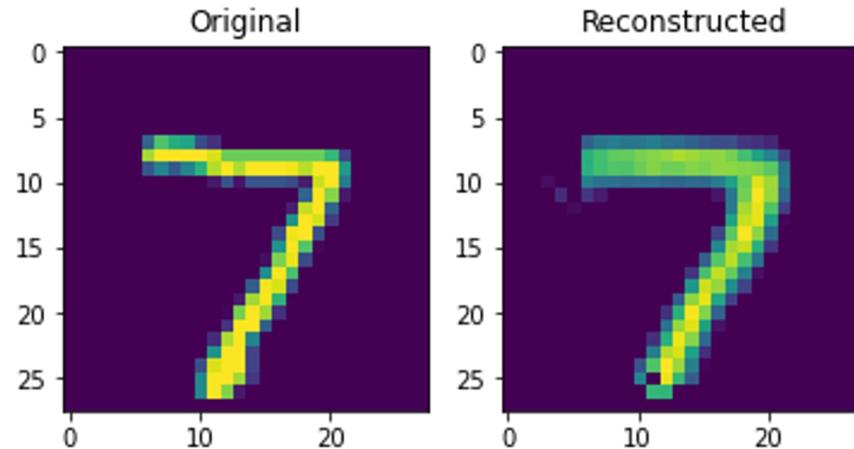
This term only depend on □

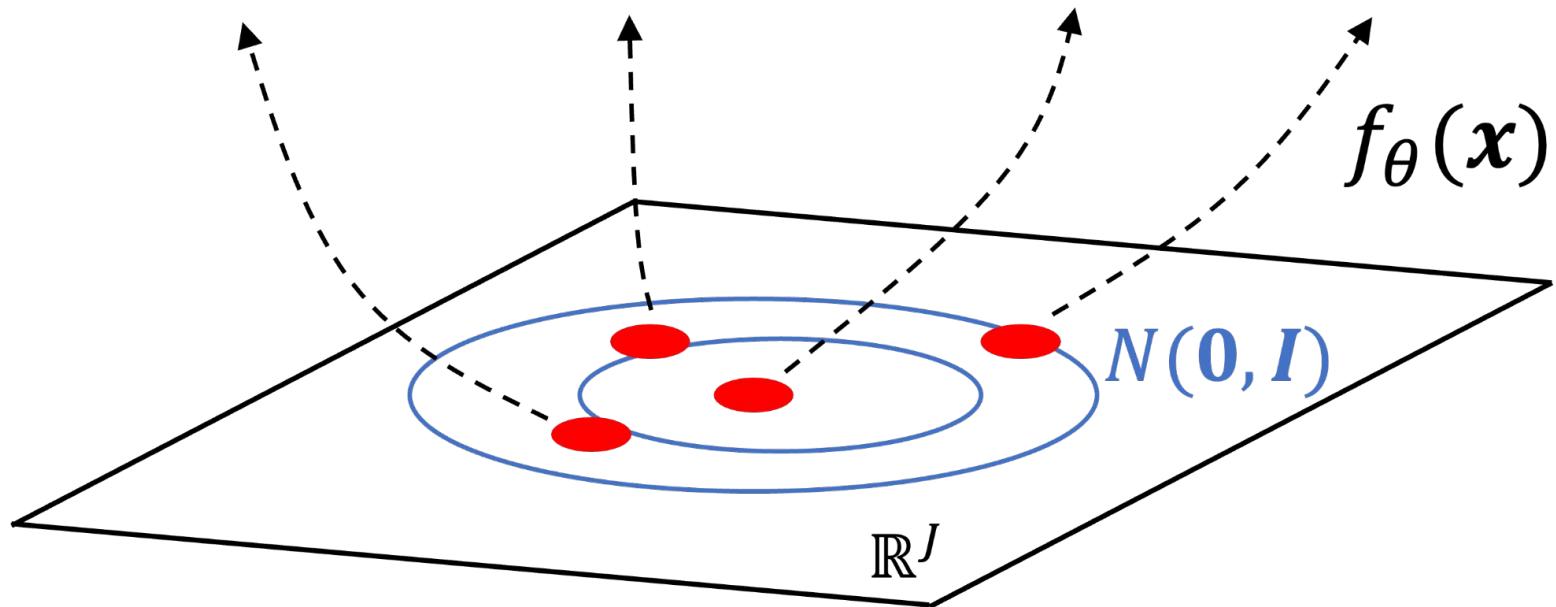
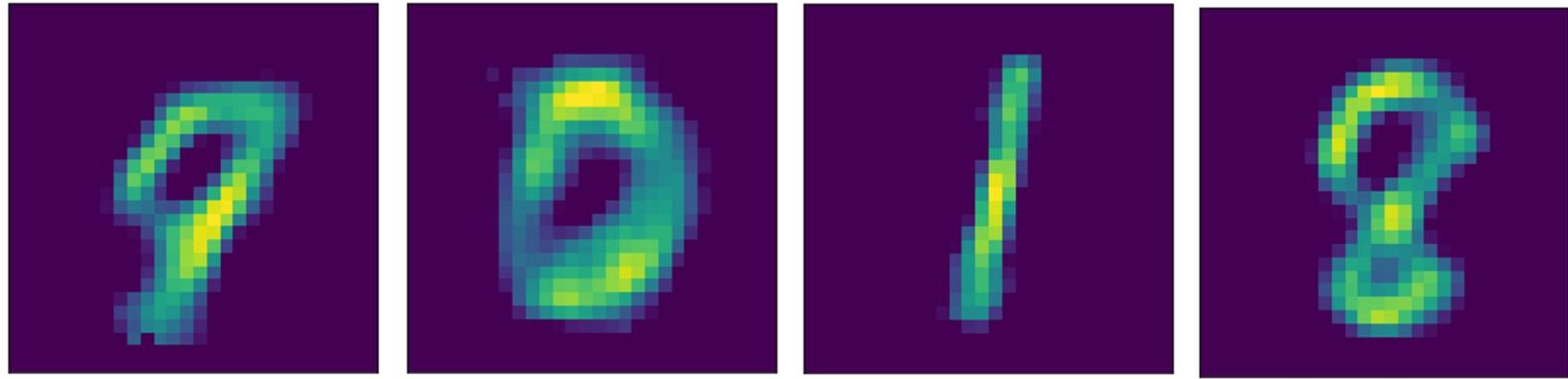
$$\max_{\theta, \varphi} \mathcal{L}_{\theta, \varphi}(x) = \underbrace{\mathbb{E}_{z \sim q_\varphi} [\log p_\theta(x|z)] - d_{KL}(q_\varphi(z|x) \| p(z))}$$

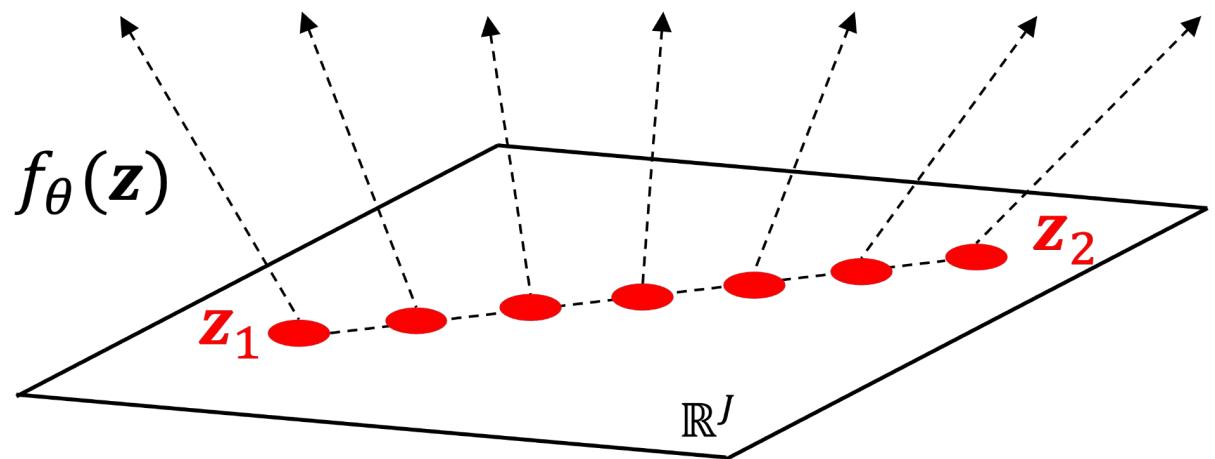
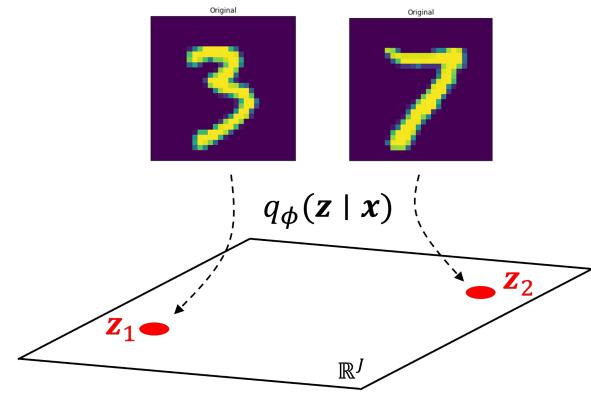
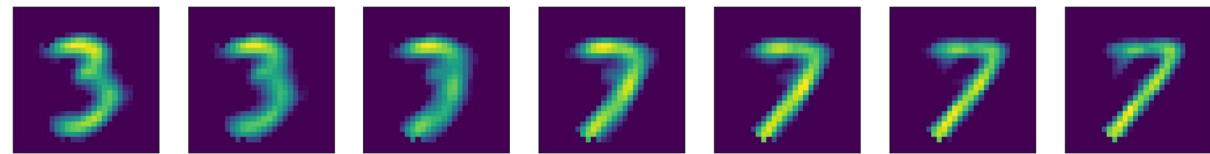
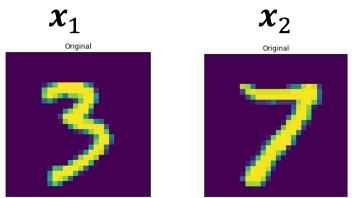
This term depends on
both θ and φ



We will introduce how we get this
in next lecture.







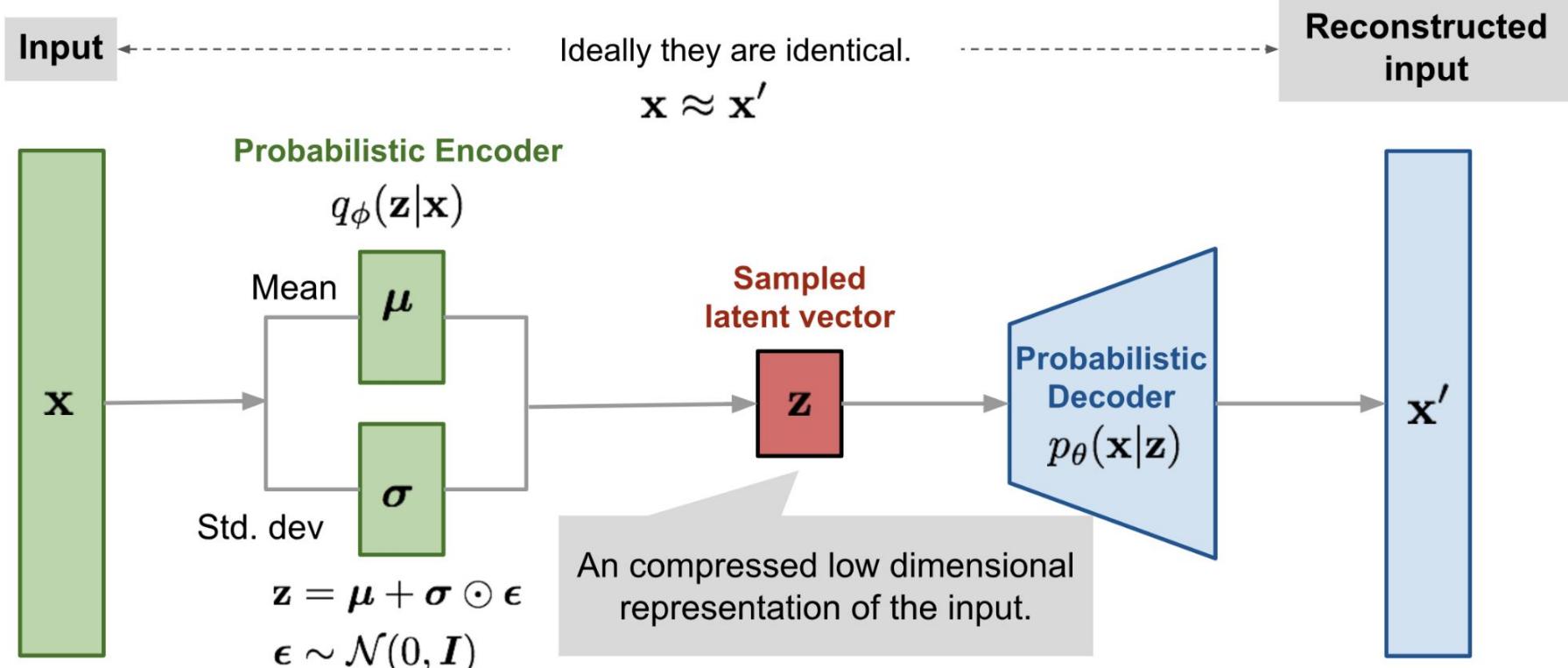
VAE Generations

The diagonal covariance constraint on the posterior forces the model to disentangle factors of variation in z .

Vary z_2 (degree of smile)



Vary z_1 (head pose)⁶⁵



We want 3 things from generative models

- *Density estimation*: Given a datapoint x , what is the probability assigned by the model, i.e. $p_\theta(x)$
- *Sampling*: How can we generate novel data from the model distribution, i.e. $x_{new} \sim p_\theta(x)$
- *Unsupervised representation learning*: How can we learn meaningful feature representations for a datapoint x ?

We want 3 things from generative models

- *Density estimation:* Given a datapoint x , what is the probability assigned by the model, i.e. $p_\theta(x)$

Roughly! You can estimate the joint $p_\theta(x, z)$ or a lower bound on $p_\theta(x)$.

Optimizing the lower bound is perhaps why VAE images are often blurrier than the ones generated with autoregressive models.

We want 3 things from generative models

- *Density estimation*: Given a datapoint x , what is the probability assigned by the model, i.e. $p_\theta(x)$
- *Sampling*: How can we generate novel data from the model distribution, i.e. $x_{new} \sim p_\theta(x)$

We can sample new data fast with one forward pass through the decoder!

We want 3 things from generative models

- *Density estimation*: Given a datapoint x , what is the probability assigned by the model, i.e. $p_\theta(x)$
- *Sampling*: How can we generate novel data from the model distribution, i.e. $x_{new} \sim p_\theta(x)$
- *Unsupervised representation learning*: How can we learn meaningful feature representations for a datapoint x ?

Yassss! The pro of VAEs is that you learn a latent code z which is produced by the encoder.

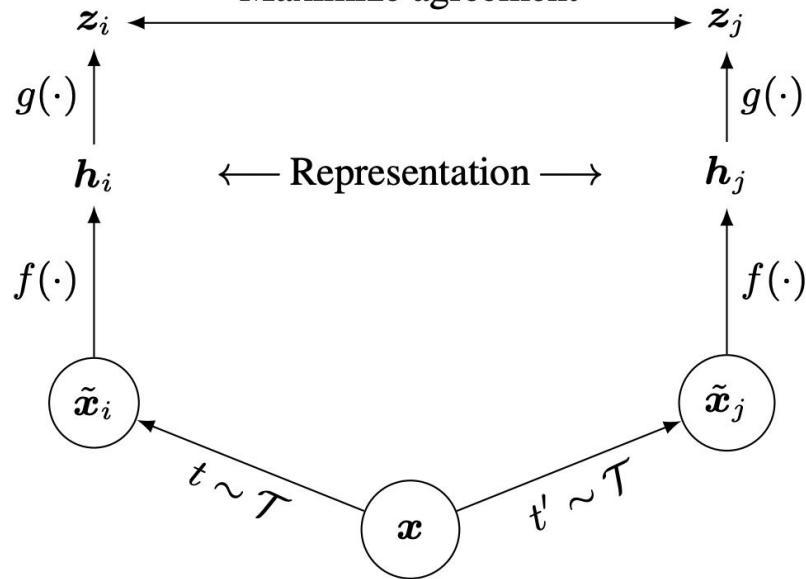
Contrastive Learning

The Contrastive Paradigm:

Train a classifier where each instance in the training set is its own class.

AKA *instance discrimination*

Maximize agreement

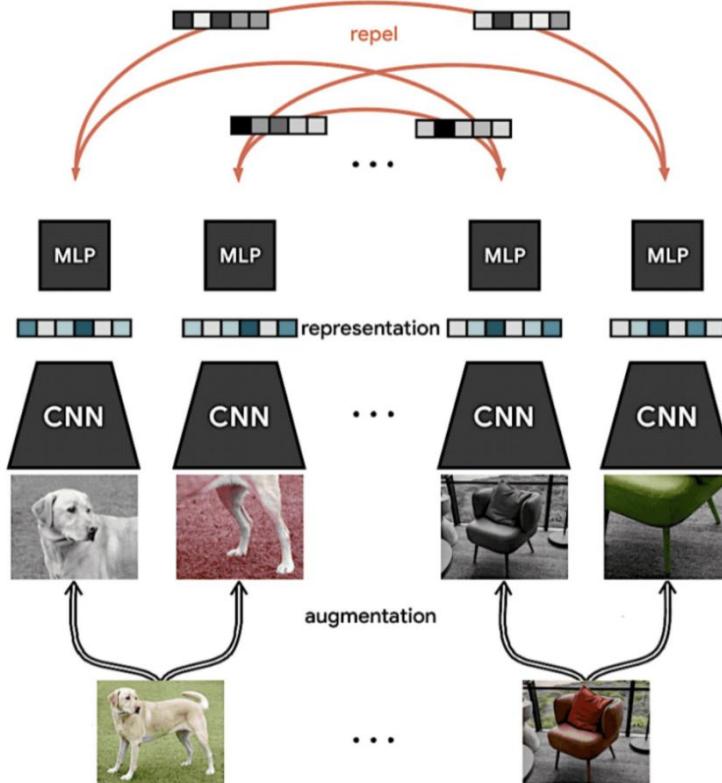


SimCLR: A Simple Framework for Contrastive Learning of Visual Representations

<https://arxiv.org/pdf/2002.05709>

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$

Modern SSL | The Contrastive Paradigm

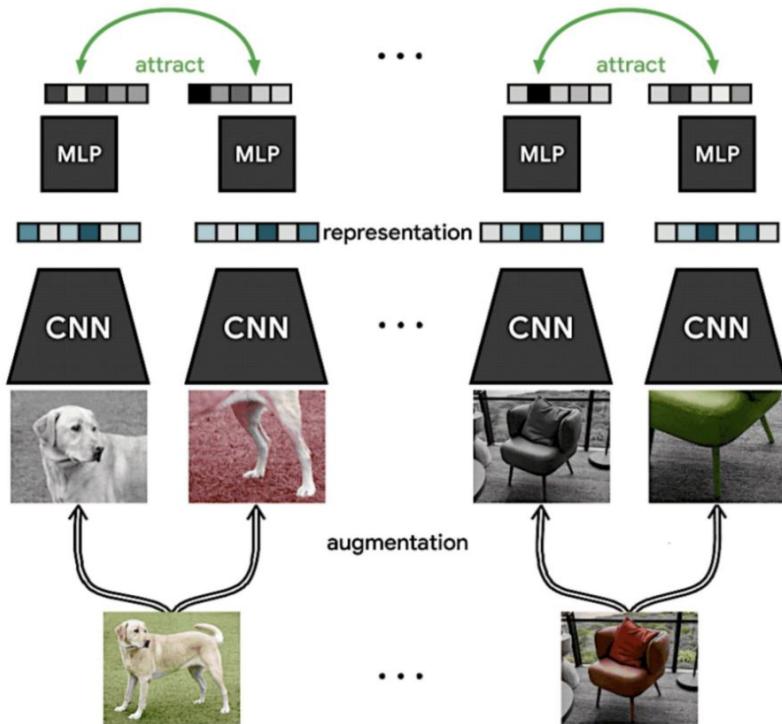


Algorithm 1 SimCLR's main learning algorithm.

input: batch size N , constant τ , structure of f, g, \mathcal{T} .
for sampled minibatch $\{x_k\}_{k=1}^N$ **do**

end for
return encoder network $f(\cdot)$, and throw away $g(\cdot)$

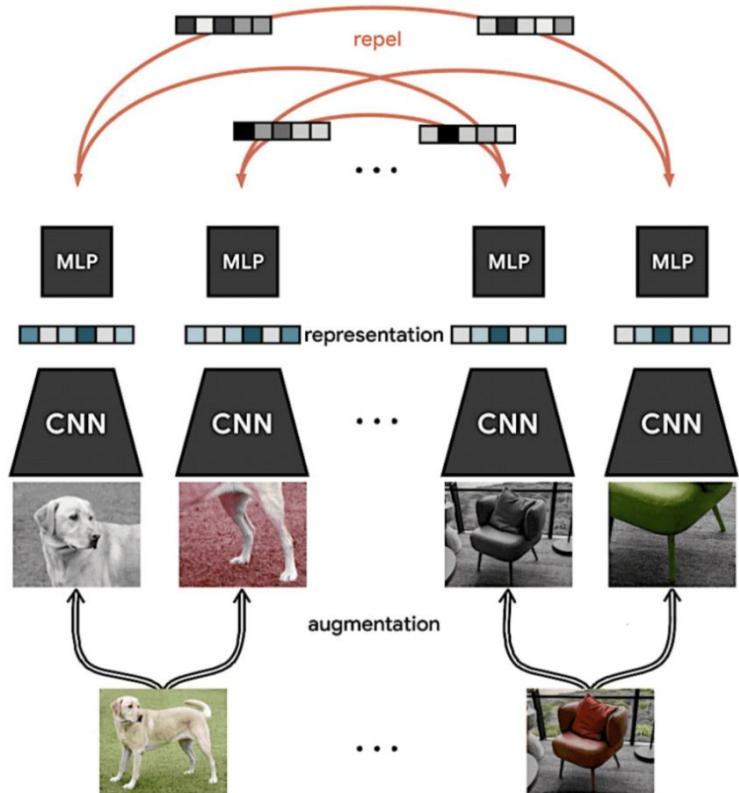
Modern SSL | The Contrastive Paradigm



Algorithm 1 SimCLR's main learning algorithm.

```
input: batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .  
for sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  do  
  for all  $k \in \{1, \dots, N\}$  do  
    draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$   
    # the first augmentation  
     $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$  # representation  
     $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # projection  
     $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$   
    # the second augmentation  
     $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$  # representation  
     $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # projection  
     $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$   
  end for  
  
end for  
return encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$ 
```

Modern SSL | The Contrastive Paradigm



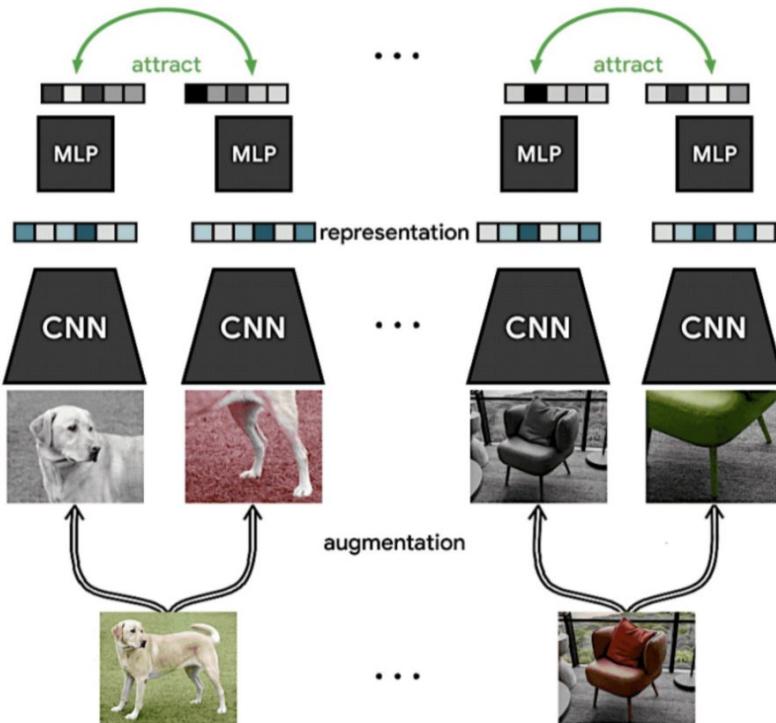
Algorithm 1 SimCLR's main learning algorithm.

```

input: batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .
for sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  do
    for all  $k \in \{1, \dots, N\}$  do
        draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$ 
        # the first augmentation
         $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$                                 # representation
         $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$                       # projection
         $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$                           # projection
        # the second augmentation
         $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$                                # representation
         $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$                          # representation
         $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$                             # projection
    end for
    for all  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  do
         $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$       # pairwise similarity
    end for
end for
return encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$ 

```

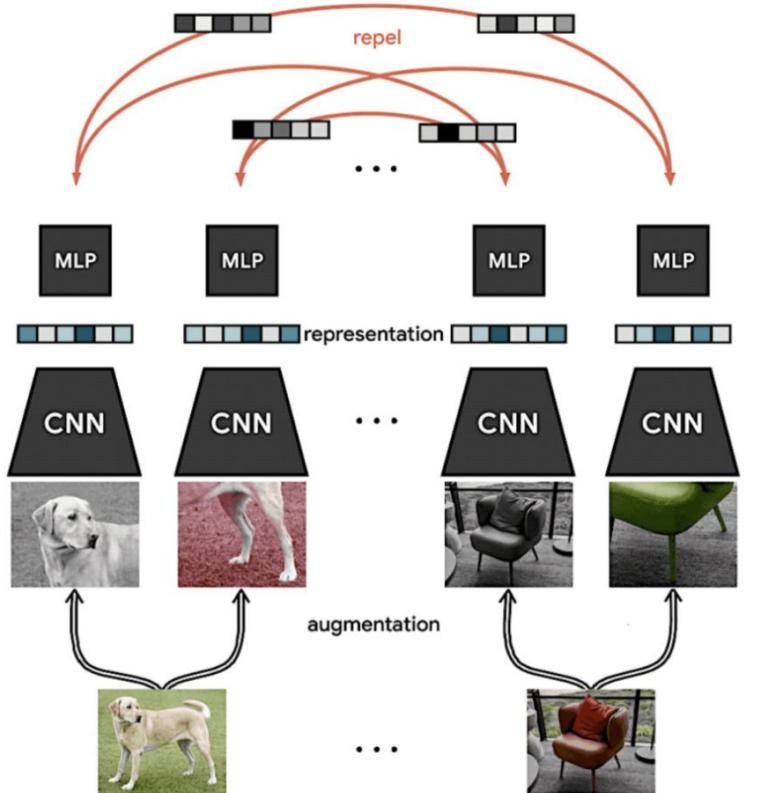
Modern SSL | The Contrastive Paradigm



Algorithm 1 SimCLR's main learning algorithm.

```
input: batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .  
for sampled minibatch  $\{x_k\}_{k=1}^N$  do  
  for all  $k \in \{1, \dots, N\}$  do  
    draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$   
    # the first augmentation  
     $\tilde{x}_{2k-1} = t(x_k)$   
     $\mathbf{h}_{2k-1} = f(\tilde{x}_{2k-1})$  # representation  
     $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection  
    # the second augmentation  
     $\tilde{x}_{2k} = t'(x_k)$   
     $\mathbf{h}_{2k} = f(\tilde{x}_{2k})$  # representation  
     $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection  
  end for  
  for all  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  do  
     $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # pairwise similarity  
  end for  
  define  $\ell(i, j)$  as  $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$   
end for  
return encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$ 
```

Modern SSL | The Contrastive Paradigm



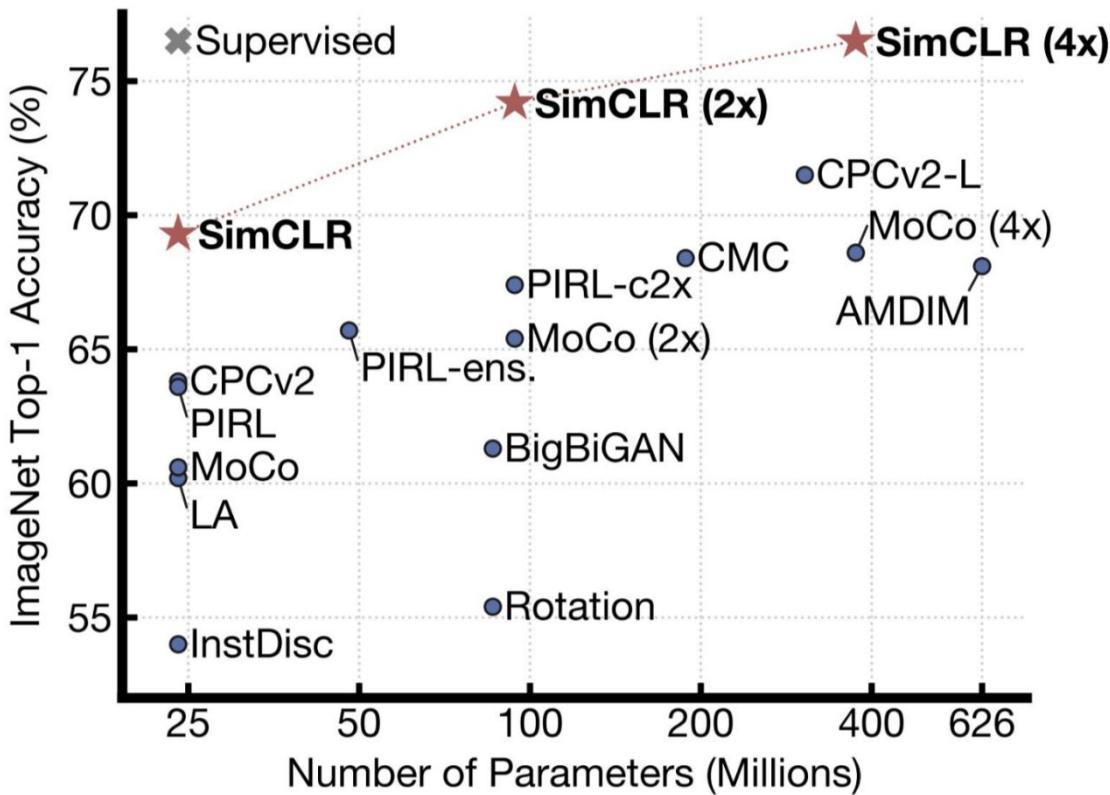
Algorithm 1 SimCLR's main learning algorithm.

```

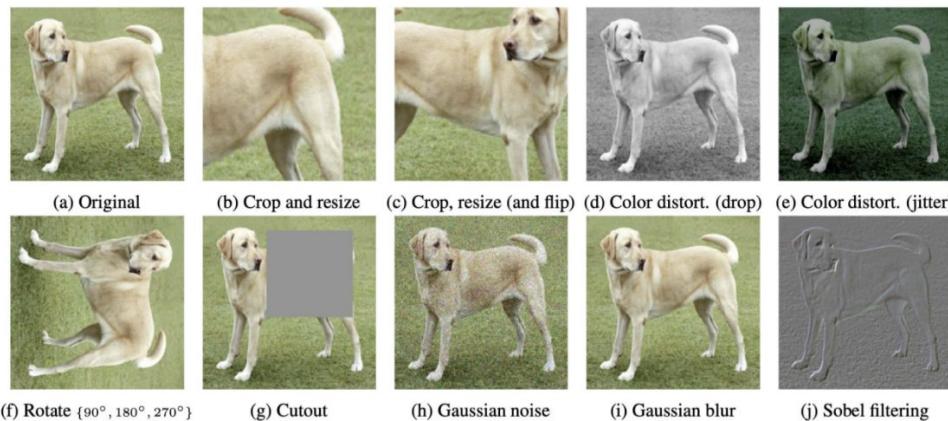
input: batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .
for sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  do
    for all  $k \in \{1, \dots, N\}$  do
        draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$ 
        # the first augmentation
         $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$ 
         $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # representation
         $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection
        # the second augmentation
         $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$ 
         $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # representation
         $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection
    end for
    for all  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  do
         $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # pairwise similarity
    end for
    define  $\ell(i, j)$  as  $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$ 
     $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$ 
    update networks  $f$  and  $g$  to minimize  $\mathcal{L}$ 
end for
return encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$ 

```

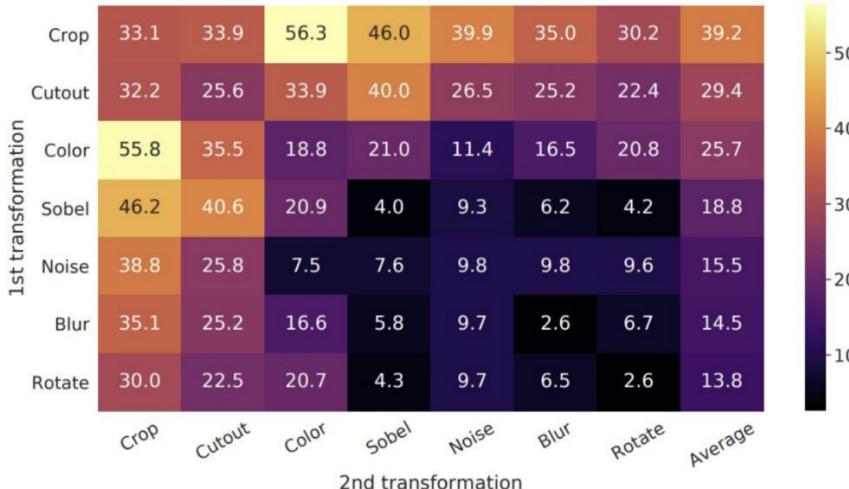
Modern SSL | The Contrastive Paradigm



Modern SSL | The Contrastive Paradigm

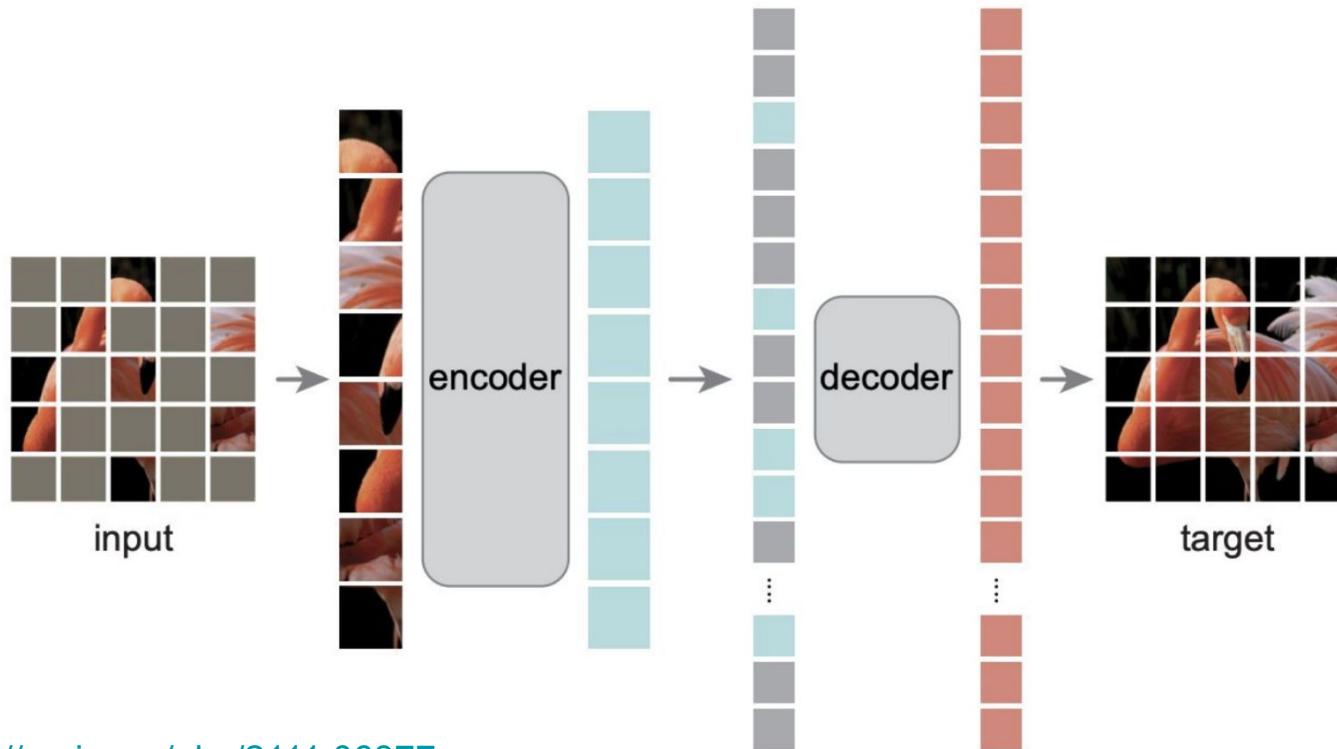


Composite
Augmentation is the
key!

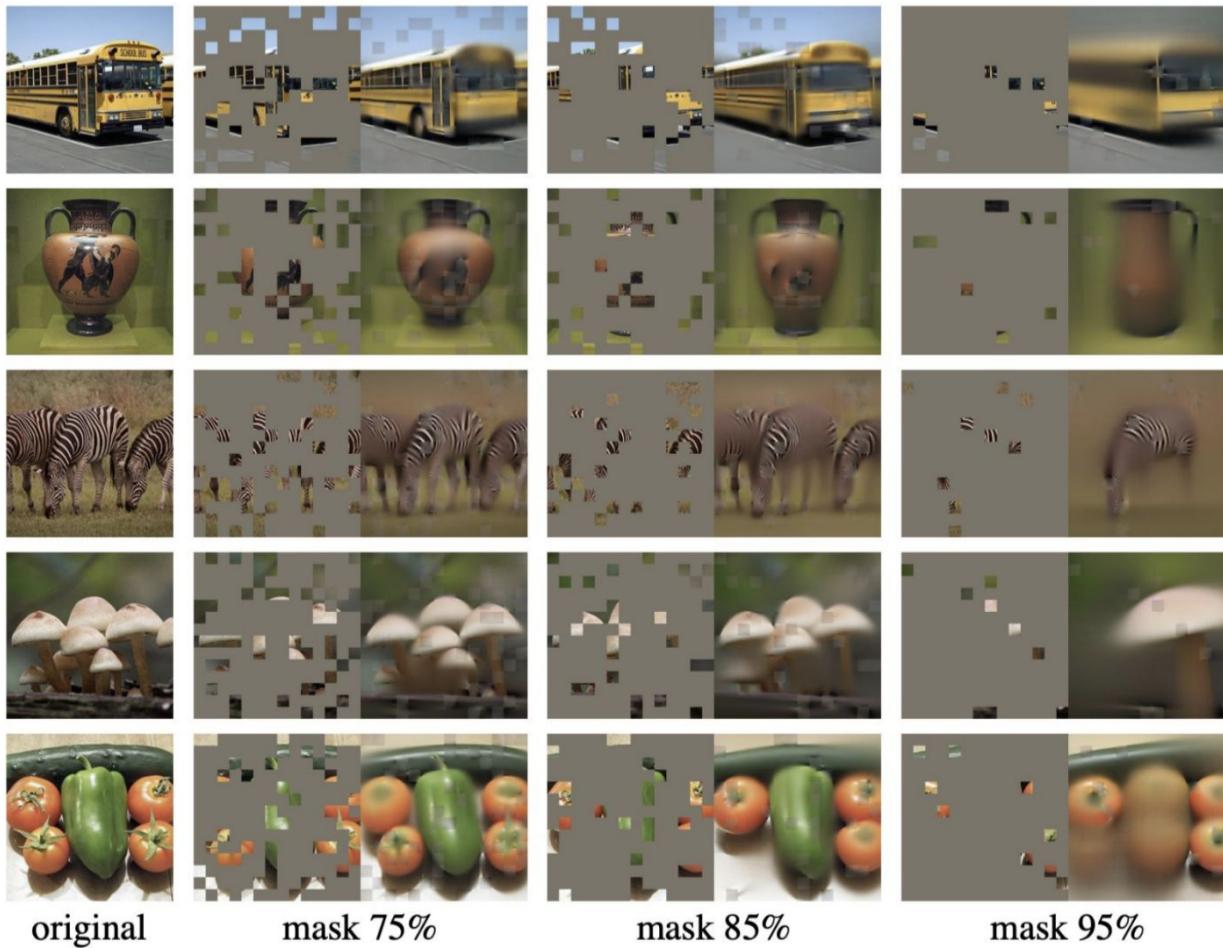


Modern SSL | The Masking Paradigm | MAE

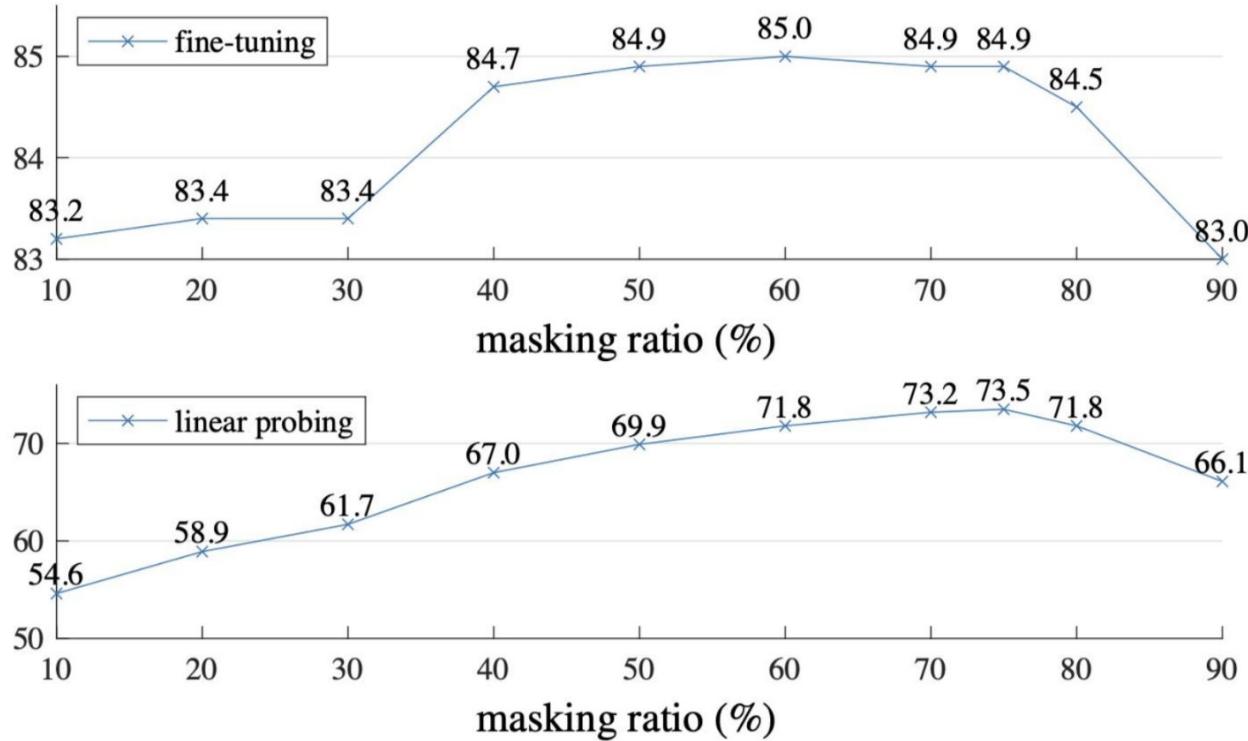
Masked Autoencoders Are Scalable Vision Learners
He et al. CVPR 2022



Modern SSL | The Masking Paradigm | MAE



Modern SSL | The Masking Paradigm | MAE



Intuition: Masking a large portion of the image to "create a challenging self supervisory task that requires holistic understanding beyond low-level image statistics."