

Variational Autoencoder from scratch

Umar Jamil

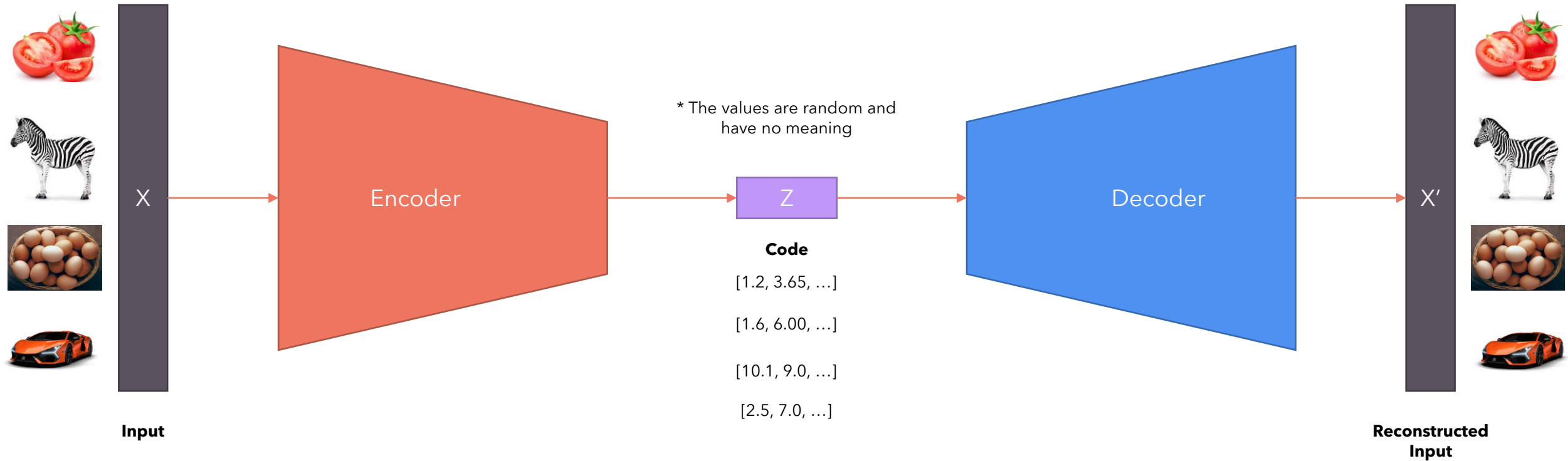
License: Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0):

<https://creativecommons.org/licenses/by-nc/4.0/legalcode>

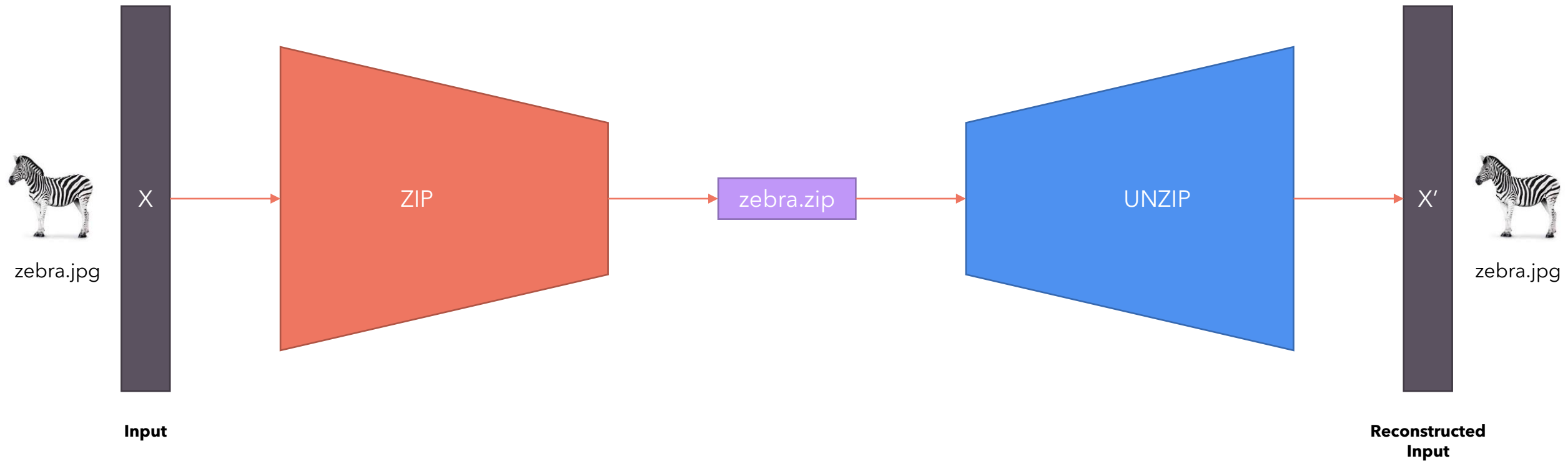
Video: <https://youtu.be/iwEzwTTalbg>

Not for commercial use

What is an Autoencoder?



Analogy with file compression

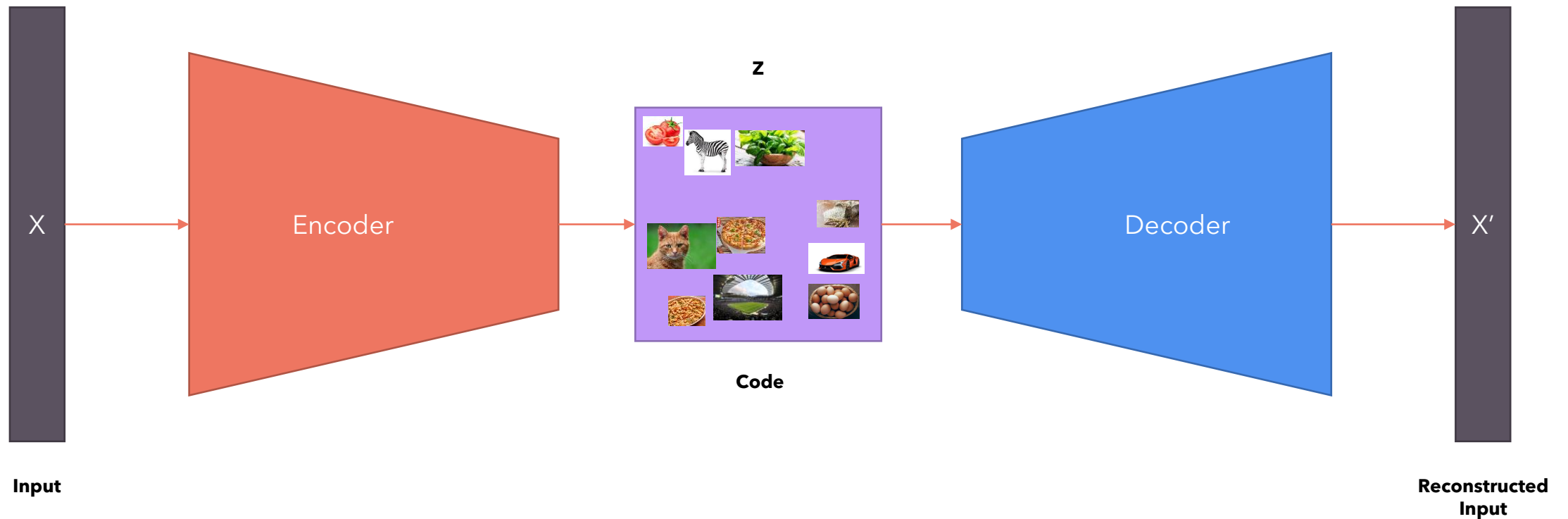


What makes a good Autoencoder?

- The code should be as small as possible, that is, the dimension of the Z vector should be as small as possible.
- The reconstructed input should as close as possible to the original input.

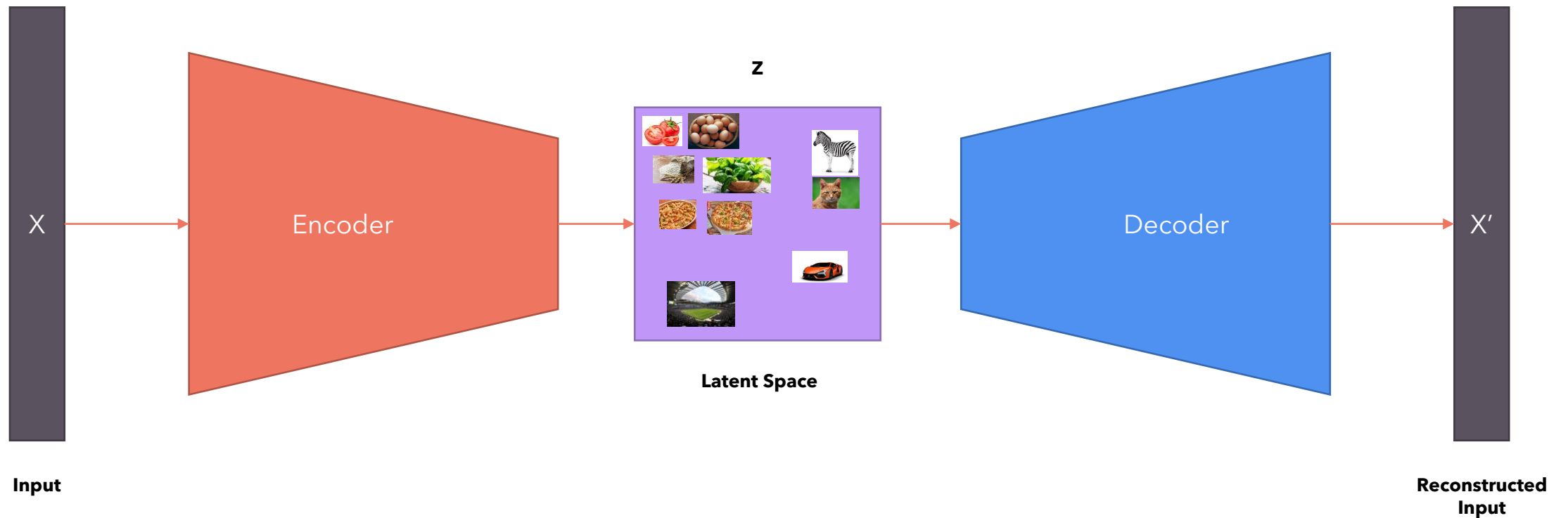
What's the problem with Autoencoders?

The code learned by the model **makes no sense**. That is, the model can just assign any vector to the inputs without the numbers in the vector representing any pattern. The model doesn't capture any **semantic relationship** between the data.



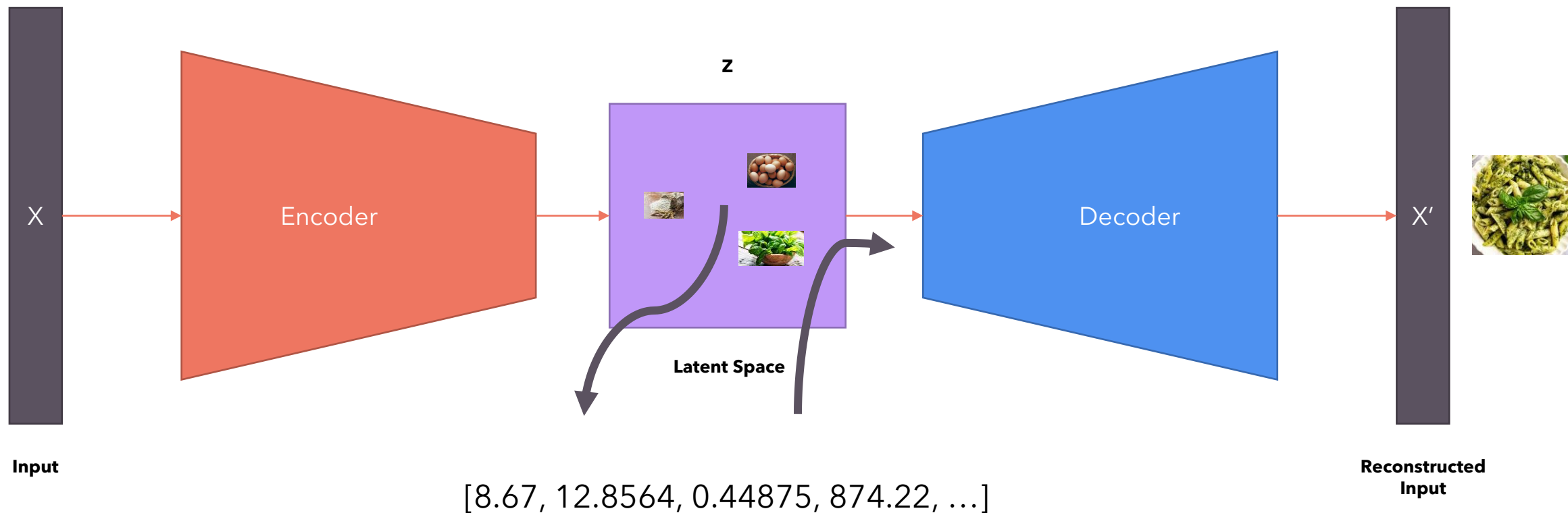
Introducing the Variational Autoencoder

The variational autoencoder, instead of learning a code, learns a "**latent space**". The latent space represents the parameters of a (multivariate) distribution.

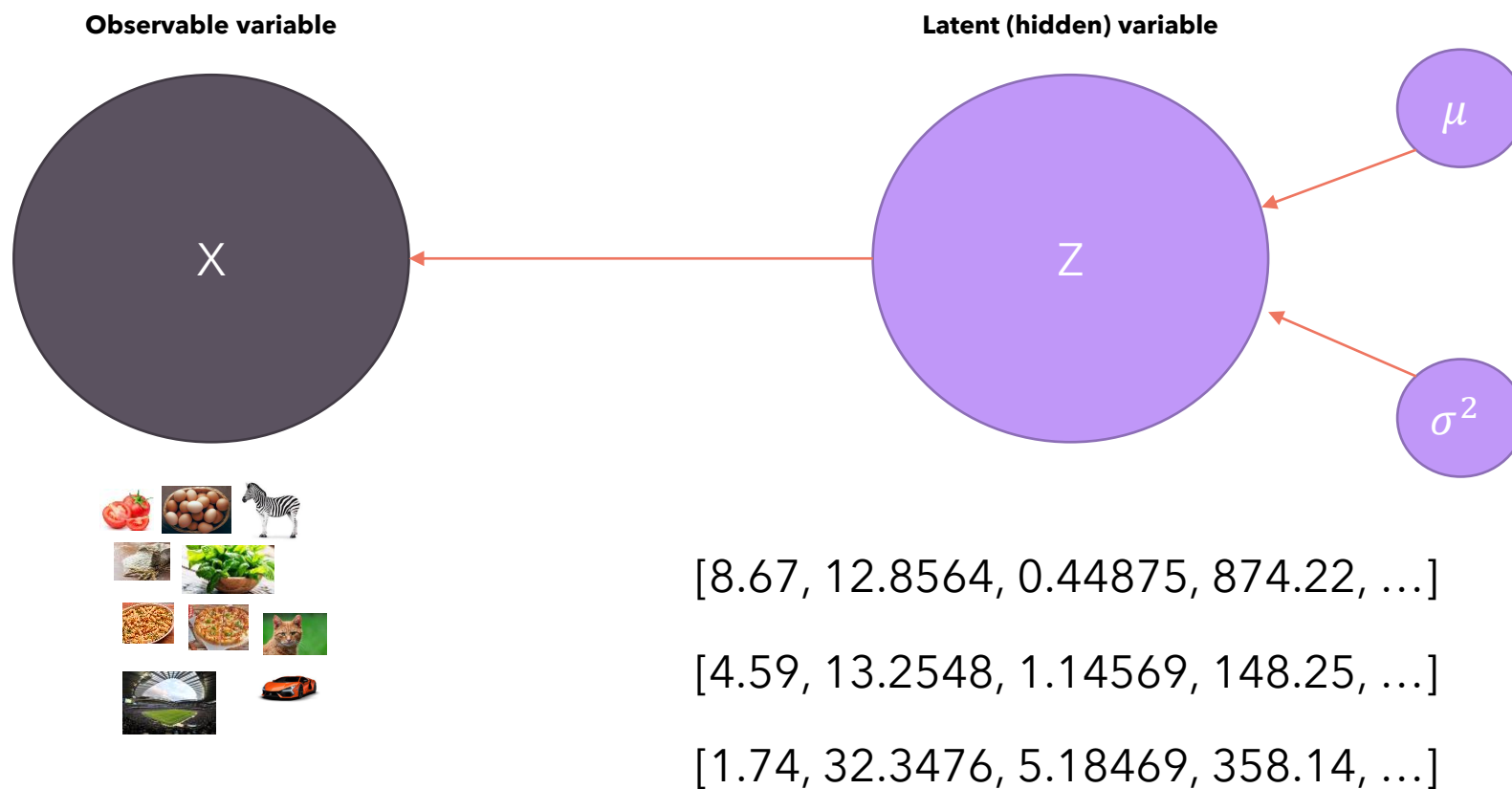


Sampling the latent space

Just like when you use Python to generate a random number between 1 and 100, you're sampling from a uniform (pseudo)random distribution between 1 and 100. In the same way, we can sample from the latent space in order to generate a random vector, give it to the decoder and generate new data.

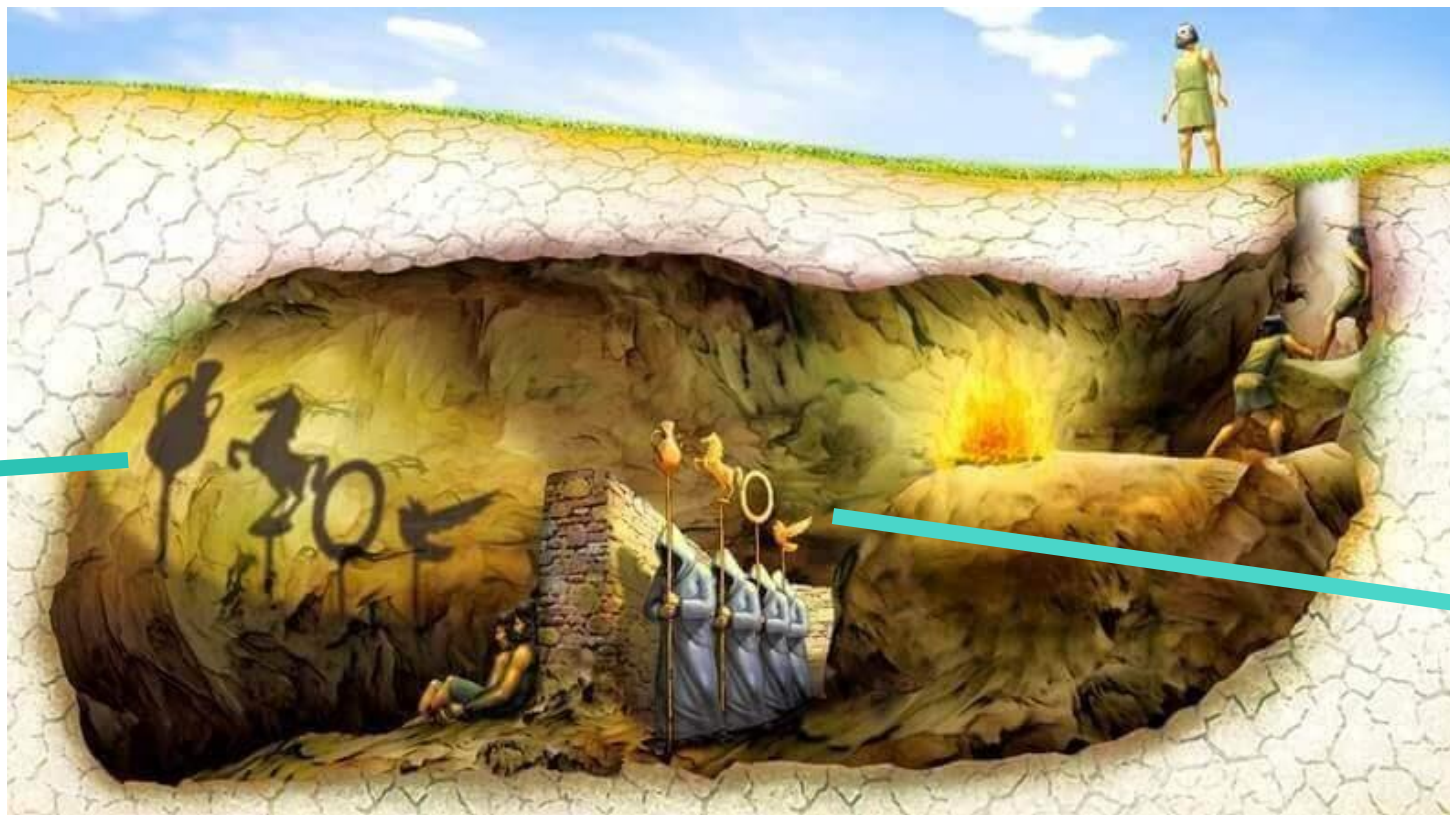
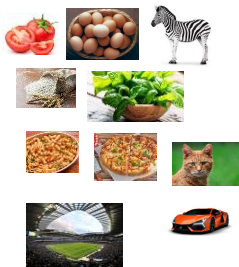


Why is it called latent space?



Plato's allegory of the cave

Observable variable



Latent (hidden) variable

[8.67, 12.8564, 0.44875, 874.22, ...]

[4.59, 13.2548, 1.14569, 148.25, ...]

[1.74, 32.3476, 5.18469, 358.14, ...]

Pep talk

1. VAE is the most important component of Stable Diffusion models. Concepts like ELBO also come in Stable Diffusion.
2. In 2023 you shouldn't be memorizing things without understanding, ChatGPT can do that faster and better than any human being. You need to be human to compete with a machine, you can't compete with a machine by acting like one.
3. You should try to learn how things work not only for curiosity, but because that's the true engine of innovation and creativity.
4. Math is fun.

Math Concepts

Expectation of a random variable

$$E_x[f(x)] = \int xf(x)dx$$

Chain rule of probability

$$P(x, y) = P(x|y)P(y)$$

Bayes' Theorem

$$P(x | y) = \frac{P(y|x)P(x)}{P(y)}$$

Kullback-Leibler Divergence

$$D_{KL}(P\|Q) = \int p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$

Properties:

- Not symmetric.
- Always ≥ 0
- It is equal to 0 if and only if $P = Q$

Let's define our model

We can define the likelihood of our data as the marginalization over the joint probability with respect to the latent variable

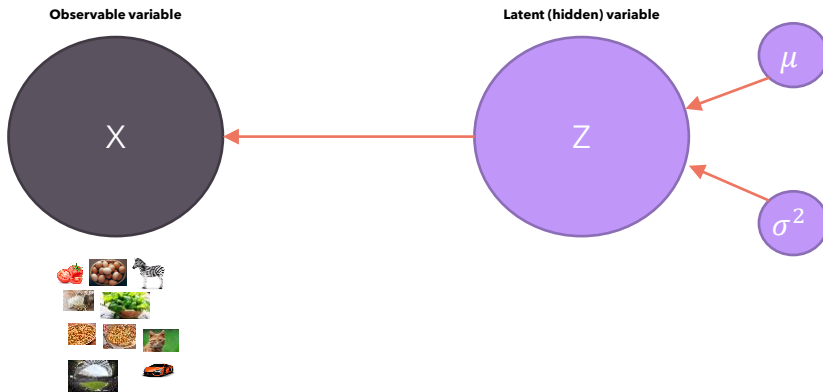
$$p(x) = \int p(x, z) dz$$

Is **intractable** because we would need to evaluate this integral over all latent variables Z .

... or we can use the Chain rule of probability

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

We **don't** have a ground truth $p(z|x)$
... which is also what we're trying to find!



Intractable problem = a problem that can be solved in theory (e.g. given large but finite resources, especially time), but for which in practice any solution takes too many resources to be useful, is known as an intractable problem.

A chicken and egg problem

In order to have a tractable $p(\mathbf{x})$ we need a tractable $p(\mathbf{z}|\mathbf{x})$

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

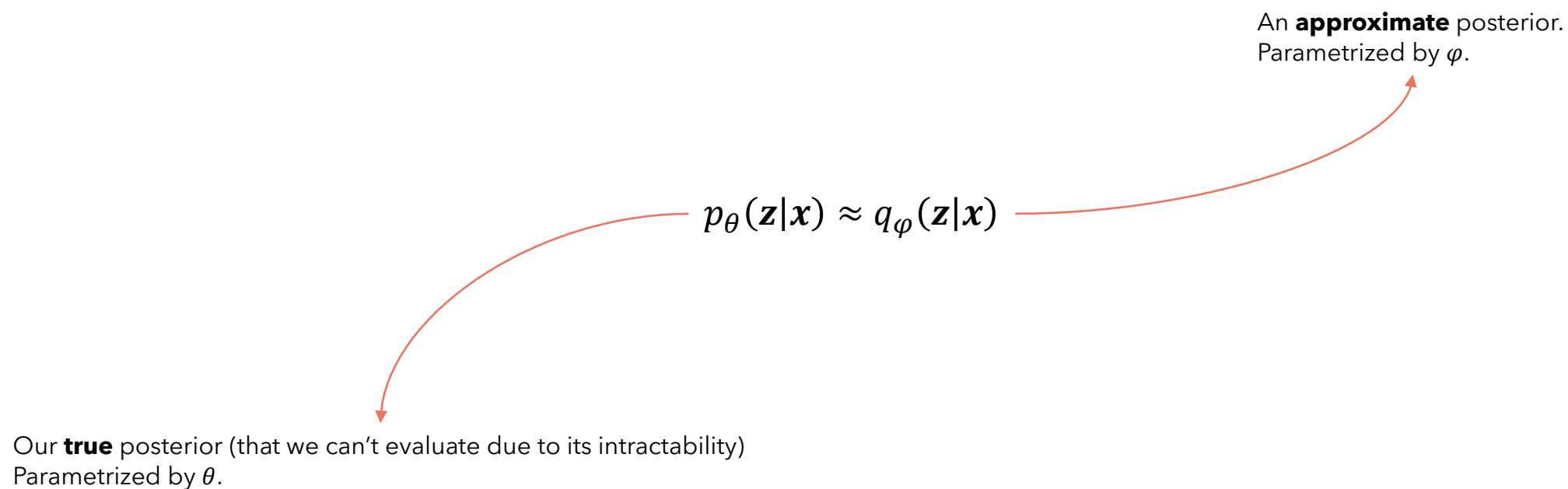


$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})}$$

In order to have a tractable $p(\mathbf{z}|\mathbf{x})$ we need a tractable $p(\mathbf{x})$



Can we find a surrogate?



Let's do some maths...

$$\log p_{\theta}(\mathbf{x}) = \log p_{\theta}(\mathbf{x})$$

$$= \log p_{\theta}(\mathbf{x}) \int q_{\varphi}(\mathbf{z}|\mathbf{x}) d\mathbf{z}$$

Multiply by 1

$$= \int \log p_{\theta}(\mathbf{x}) q_{\varphi}(\mathbf{z}|\mathbf{x}) d\mathbf{z}$$

Bring inside the integral

$$= E_{q_{\varphi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x})]$$

Definition of expectation

$$= E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right]$$

Apply the equation $p_{\theta}(\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})}$

$$= E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z}) q_{\varphi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}|\mathbf{x}) q_{\varphi}(\mathbf{z}|\mathbf{x})} \right]$$

Multiply by 1

$$= E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right] + E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_{\varphi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right]$$

Split the expectation

$$= E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right] + \underbrace{D_{KL}(q_{\varphi}(\mathbf{z}|\mathbf{x}) \| p_{\theta}(\mathbf{z}|\mathbf{x}))}_{\geq 0}$$

Definition of KL divergence

What can we infer?

$$\log p_{\theta}(\mathbf{x}) = \underbrace{E_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]}_{\text{ELBO}} + \underbrace{D_{KL} \left(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{\theta}(\mathbf{z}|\mathbf{x}) \right)}_{\geq 0}$$

ELBO = Evidence Lower Bound

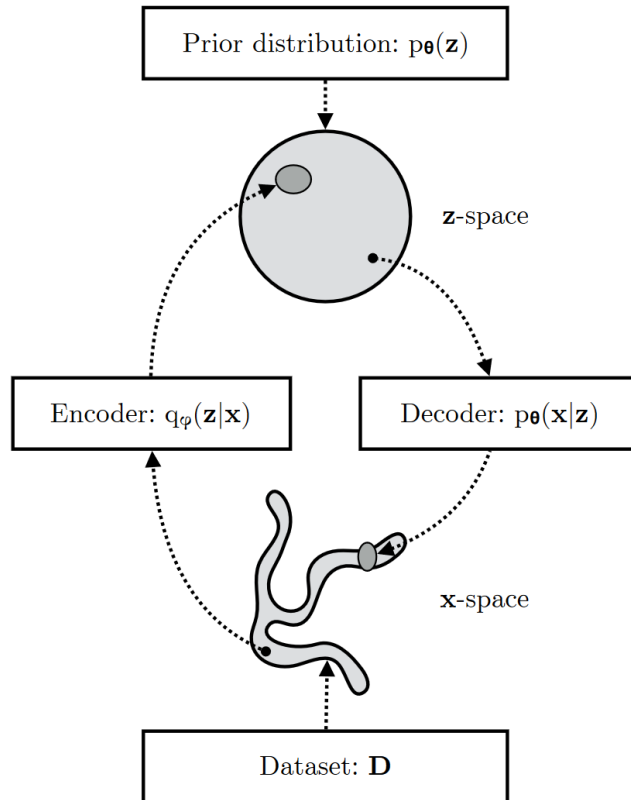
$$\text{Total Compensation} = \text{Base Salary} + \underbrace{\text{Bonus}}_{\geq 0}$$

We can for sure deduce the following:

$$\text{Total Compensation} \geq \text{Base Salary}$$

$$\log p_{\theta}(\mathbf{x}) \geq E_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

ELBO in detail



$$\log p_\theta(\mathbf{x}) \geq E_{q_\varphi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\varphi(\mathbf{z}|\mathbf{x})} \right] = E_{q_\varphi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\varphi(\mathbf{z}|\mathbf{x})} \right]$$

Chain rule of probability

$$= E_{q_\varphi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] + E_{q_\varphi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{z})}{q_\varphi(\mathbf{z}|\mathbf{x})} \right]$$

Split the expectation

$$= E_{q_\varphi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\varphi(\mathbf{z}|\mathbf{x}) \| p_\theta(\mathbf{z}))$$

Definition of KL divergence

$$\underbrace{\phantom{E_{q_\varphi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}}_{\text{Profit}} = \underbrace{\phantom{E_{q_\varphi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}}_{\text{Revenue}} - \underbrace{\phantom{D_{KL}(q_\varphi(\mathbf{z}|\mathbf{x}) \| p_\theta(\mathbf{z}))}}_{\text{Costs}}$$

Maximizing the ELBO means:

1. Maximizing the first term: maximizing the reconstruction likelihood of the decoder
2. Minimizing the second term: minimizing the distance between the learned distribution and the prior belief we have over the latent variable.

Maximizing the ELBO: A little introduction to estimators

- When we have a function we want to **maximize**, we usually take the gradient and adjust the weights of the model so that they move **along** the gradient direction.
- When we have a function we want to **minimize**, we usually take the gradient, and adjust the weights of the model so that they move **against** the gradient direction.

Stochastic Gradient Descent

When used to minimize the above function, a standard (or "batch") **gradient descent** method would perform the following iterations:

$$w := w - \eta \nabla Q(w) = w - \frac{\eta}{n} \sum_{i=1}^n \nabla Q_i(w),$$

where η is a step size (sometimes called the **learning rate** in machine learning).

SGD is **stochastic** because we choose the minibatch randomly from our dataset and we then average the loss over the minibatch

$$L(\theta, \varphi, \mathbf{x}) = E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right] = E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D_{KL} \left(q_{\varphi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}) \right)$$

How to maximize the ELBO?

ELBO

$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = -D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})) + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})] \quad (3)$$

We want to differentiate and optimize the lower bound $\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$ w.r.t. both the variational parameters ϕ and generative parameters θ . However, the gradient of the lower bound w.r.t. ϕ is a bit problematic. The usual (naïve) Monte Carlo gradient estimator for this type of problem is: $\nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z})} [f(\mathbf{z})] = \mathbb{E}_{q_{\phi}(\mathbf{z})} [f(\mathbf{z}) \nabla_{q_{\phi}(\mathbf{z})} \log q_{\phi}(\mathbf{z})] \simeq \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}) \nabla_{q_{\phi}(\mathbf{z}^{(l)})} \log q_{\phi}(\mathbf{z}^{(l)})$ where $\mathbf{z}^{(l)} \sim q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$. This gradient estimator exhibits very high variance (see e.g. [BJP12]) and is impractical for our purposes.

SCORE estimator

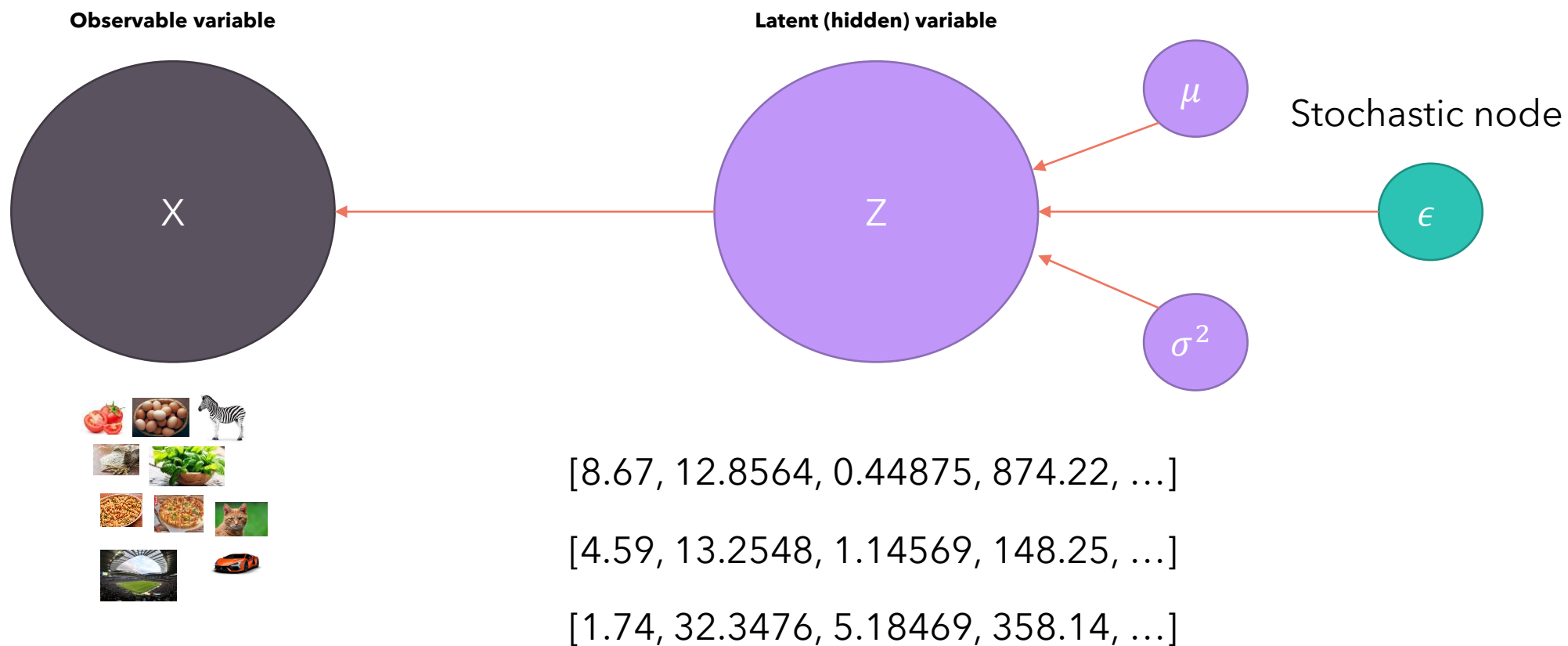
Kingma, D.P. and Welling, M., 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

This estimator is **unbiased**, meaning that even if at every step it may not be equal to the true expectation, on average it will converge to it, but as it is stochastic, it also has a variance and it happens to be high for practical use. Plus, we can't run backpropagation through it!

We need a new estimator!



The reparameterization trick



Running backpropagation on the reparametrized model

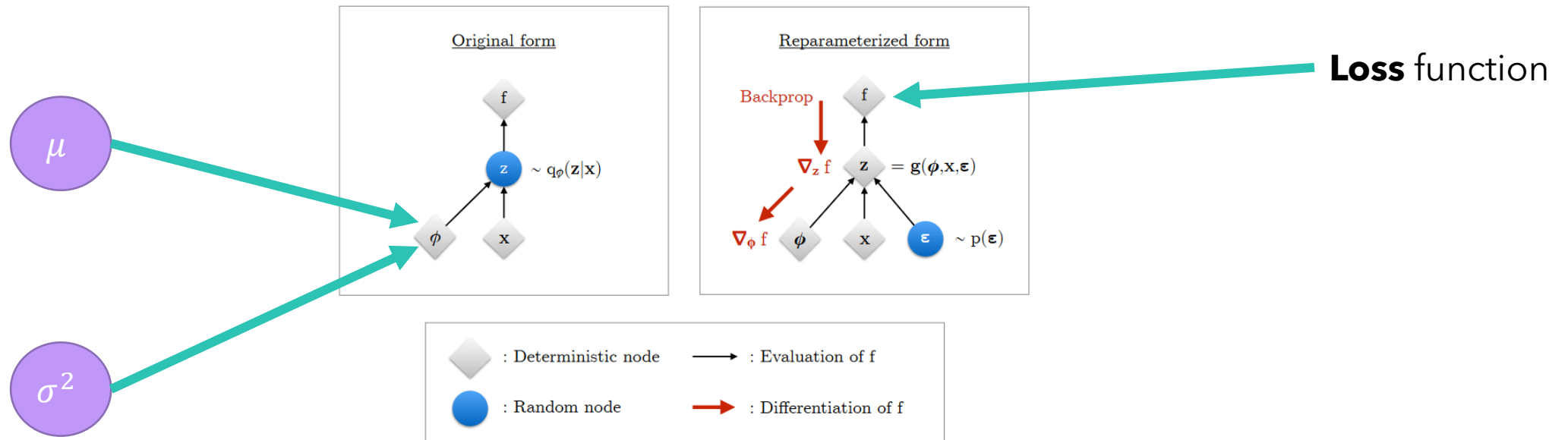


Figure 2.3: Illustration of the reparameterization trick. The variational parameters ϕ affect the objective f through the random variable $z \sim q_\phi(z|x)$. We wish to compute gradients $\nabla_\phi f$ to optimize the objective with SGD. In the original form (left), we cannot differentiate f w.r.t. ϕ , because we cannot directly backpropagate gradients through the random variable z . We can 'externalize' the randomness in z by re-parameterizing the variable as a deterministic and differentiable function of ϕ , x , and a newly introduced random variable ϵ . This allows us to 'backprop through z ', and compute gradients $\nabla_\phi f$.

Kingma, D.P. and Welling, M., 2019. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4), pp.307-392.

A new estimator!

$$L(\theta, \varphi, \mathbf{x}) = E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right] = E_{p(\epsilon)} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right]$$

$\underbrace{\hspace{10em}}_{\text{ELBO}}$

$$\epsilon \approx p(\epsilon)$$

$$\mathbf{z} = g(\varphi, \mathbf{x}, \epsilon)$$

$$E_{p(\epsilon)} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right] \cong \tilde{\mathcal{L}}(\theta, \varphi, \mathbf{x}) = \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})}$$

aging the approximate posterior to be close to the prior $p_{\theta}(\mathbf{z})$. This yields a second version of the SGVB estimator $\tilde{\mathcal{L}}^B(\theta, \phi; \mathbf{x}^{(i)}) \simeq \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$, corresponding to eq. (3), which typically has less variance than the generic estimator:

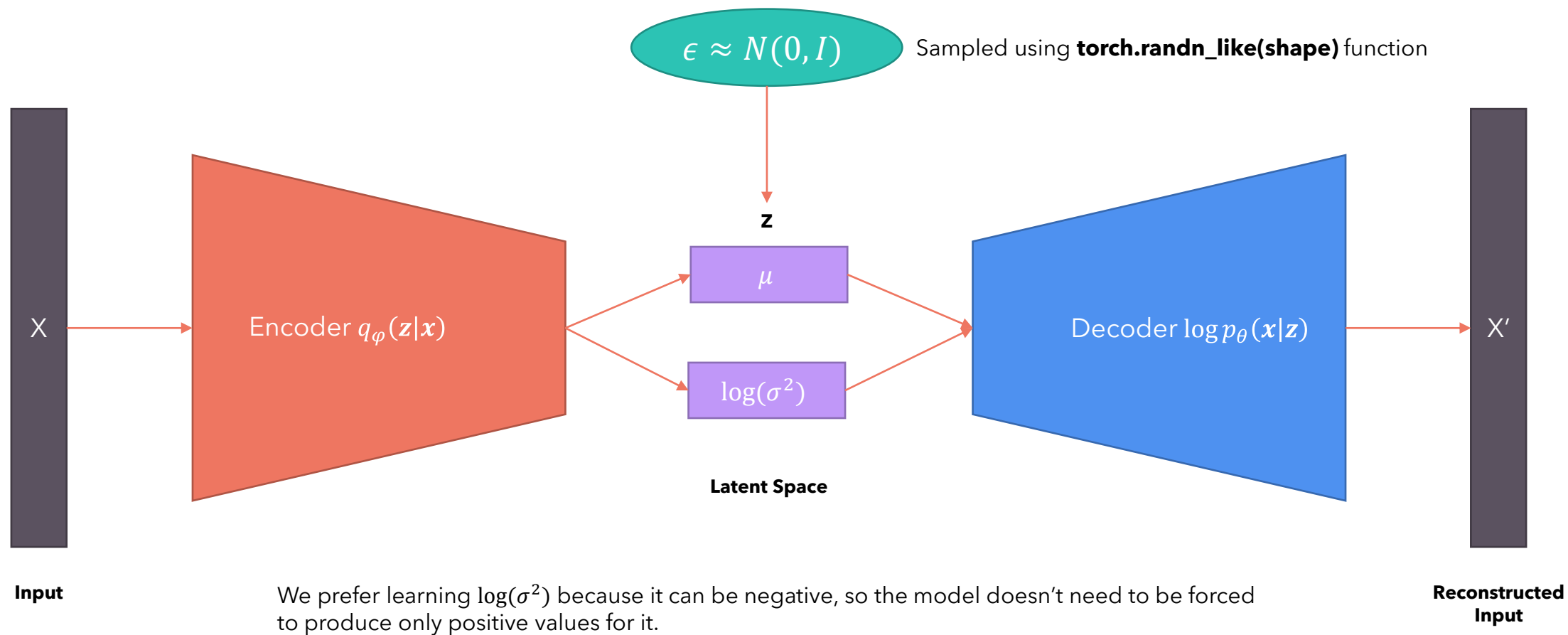
Kingma, D.P. and Welling, M., 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

Is the new estimator unbiased?

$$L(\theta, \varphi, \mathbf{x}) = E_{p(\epsilon)} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right] \quad \tilde{L}(\theta, \varphi, \mathbf{x}) = \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})}$$

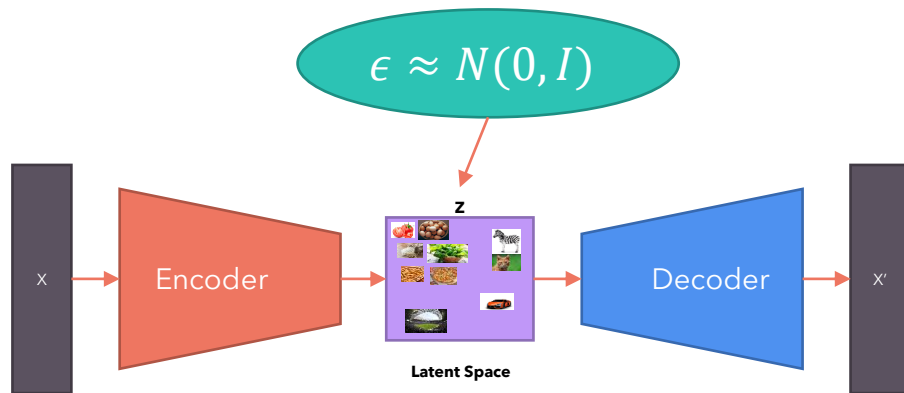
$$\begin{aligned} E_{p(\epsilon)} [\nabla_{\theta, \varphi} \tilde{L}(\theta, \varphi, \mathbf{x})] &= E_{p(\epsilon)} \left[\nabla_{\theta, \varphi} \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right] \\ &= \nabla_{\theta, \varphi} (E_{p(\epsilon)} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right]) \\ &= \nabla_{\theta, \varphi} (L(\theta, \varphi, \mathbf{x})) \end{aligned}$$

Example network



Show me the loss already!

MLP = Multi Layer Perceptron



Let the prior over the latent variables be the centered isotropic multivariate Gaussian $p_{\theta}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$. Note that in this case, the prior lacks parameters. We let $p_{\theta}(\mathbf{x}|\mathbf{z})$ be a multivariate Gaussian (in case of real-valued data) or Bernoulli (in case of binary data) whose distribution parameters are computed from \mathbf{z} with a MLP (a fully-connected neural network with a single hidden layer, see appendix C). Note the true posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$ is in this case intractable. While there is much freedom in the form $q_{\phi}(\mathbf{z}|\mathbf{x})$, we'll assume the true (but intractable) posterior takes on an approximate Gaussian form with an approximately diagonal covariance. In this case, we can let the variational approximate posterior be a multivariate Gaussian with a diagonal covariance structure:

$$\log q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)}) = \log \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{2(i)} \mathbf{I}) \quad (9)$$

where the mean and s.d. of the approximate posterior, $\boldsymbol{\mu}^{(i)}$ and $\boldsymbol{\sigma}^{(i)}$, are outputs of the encoding MLP, i.e. nonlinear functions of datapoint $\mathbf{x}^{(i)}$ and the variational parameters ϕ (see appendix C).

As explained in section 2.4, we sample from the posterior $\mathbf{z}^{(i,l)} \sim q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$ using $\mathbf{z}^{(i,l)} = g_{\phi}(\mathbf{x}^{(i)}, \boldsymbol{\epsilon}^{(l)}) = \boldsymbol{\mu}^{(i)} + \boldsymbol{\sigma}^{(i)} \odot \boldsymbol{\epsilon}^{(l)}$ where $\boldsymbol{\epsilon}^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. With \odot we signify an element-wise product. In this model both $p_{\theta}(\mathbf{z})$ (the prior) and $q_{\phi}(\mathbf{z}|\mathbf{x})$ are Gaussian; in this case, we can use the estimator of eq. (7) where the KL divergence can be computed and differentiated without estimation (see appendix B). The resulting estimator for this model and datapoint $\mathbf{x}^{(i)}$ is:

$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) \simeq \frac{1}{2} \sum_{j=1}^J \left(1 + \log((\sigma_j^{(i)})^2) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2 \right) + \frac{1}{L} \sum_{l=1}^L \boxed{\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}^{(i,l)})} \quad (10)$$

where $\mathbf{z}^{(i,l)} = \boldsymbol{\mu}^{(i)} + \boldsymbol{\sigma}^{(i)} \odot \boldsymbol{\epsilon}^{(l)}$ and $\boldsymbol{\epsilon}^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

As explained above and in appendix C, the decoding term $\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}^{(i,l)})$ is a Bernoulli or Gaussian MLP, depending on the type of data we are modelling.

Kingma, D.P. and Welling, M., 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

How to derive the loss function?



48



The encoder distribution is $q(z|x) = \mathcal{N}(z|\mu(x), \Sigma(x))$ where $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$.

The latent prior is given by $p(z) = \mathcal{N}(0, I)$.

Both are multivariate Gaussians of dimension n , for which in general the KL divergence is:

$$\mathcal{D}_{\text{KL}}[p_1 || p_2] = \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - n + \text{tr}\{\Sigma_2^{-1} \Sigma_1\} + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right]$$

where $p_1 = \mathcal{N}(\mu_1, \Sigma_1)$ and $p_2 = \mathcal{N}(\mu_2, \Sigma_2)$.

In the VAE case, $p_1 = q(z|x)$ and $p_2 = p(z)$, so $\mu_1 = \mu$, $\Sigma_1 = \Sigma$, $\mu_2 = \vec{0}$, $\Sigma_2 = I$. Thus:

$$\begin{aligned} \mathcal{D}_{\text{KL}}[q(z|x) || p(z)] &= \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - n + \text{tr}\{\Sigma_2^{-1} \Sigma_1\} + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right] \\ &= \frac{1}{2} \left[\log \frac{|I|}{|\Sigma|} - n + \text{tr}\{I^{-1} \Sigma\} + (\vec{0} - \mu)^T I^{-1} (\vec{0} - \mu) \right] \\ &= \frac{1}{2} \left[-\log |\Sigma| - n + \text{tr}\{\Sigma\} + \mu^T \mu \right] \\ &= \frac{1}{2} \left[-\log \prod_i \sigma_i^2 - n + \sum_i \sigma_i^2 + \sum_i \mu_i^2 \right] \\ &= \frac{1}{2} \left[-\sum_i \log \sigma_i^2 - n + \sum_i \sigma_i^2 + \sum_i \mu_i^2 \right] \\ &= \frac{1}{2} \left[-\sum_i (\log \sigma_i^2 + 1) + \sum_i \sigma_i^2 + \sum_i \mu_i^2 \right] \end{aligned}$$

<https://stats.stackexchange.com/questions/318748/deriving-the-kl-divergence-loss-for-vaes>

Thanks for watching!
Don't forget to subscribe for
more amazing content on AI
and Machine Learning!