

A brief introduction to (variational) autoencoders

Alexej Klushyn

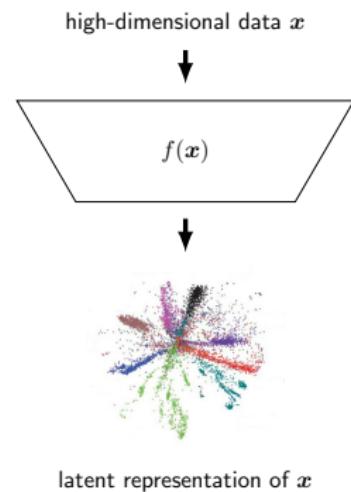
argmax.ai

Volkswagen Group Machine Learning Research Lab

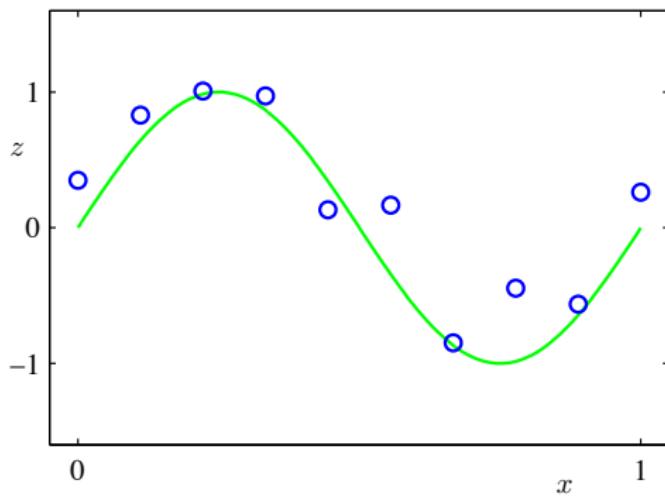
Can we learn a latent representation of our data?

- ▶ Dimensionality reduction methods map high-dimensional data to a low-dimensional latent space.
- ▶ Often this mapping is linear or restricted in some other way.

Can we learn an arbitrary complex mapping $f(x)$?



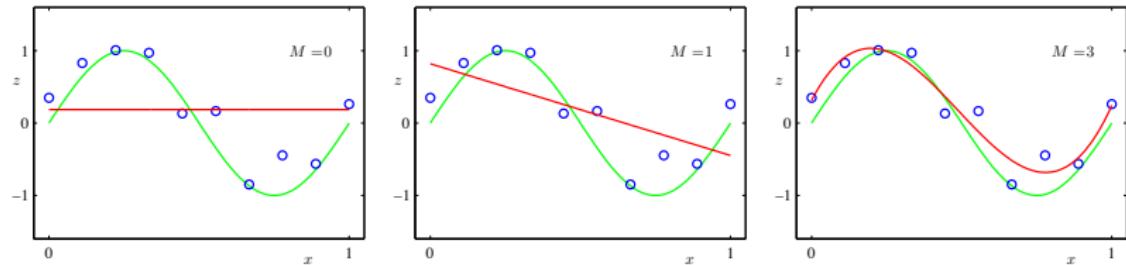
A noisy real-valued function



inputs: $\mathbf{x} = (x_1, \dots, x_N)^T$

targets: $\mathbf{z} = (z_1, \dots, z_N)^T, \quad z_i = h(x_i) + \epsilon = \sin(2\pi x_i) + \epsilon$

Model: 0th, 1st, and 3rd order polynomial



0th order: $y(\mathbf{x}, \mathbf{w}) = w_0$

1st order: $y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x$

3rd order: $y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$

Problem definition

We have input vectors \mathbf{x} and associated output vectors \mathbf{z} . We want to describe the underlying functional relation.

What about the following model?

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

where

- | | | | |
|--------|----------------|---|--|
| ϕ | basis function | — | many choices, can be nonlinear |
| w_0 | bias | — | equivalent to defining $\phi_0 \equiv 1$ |

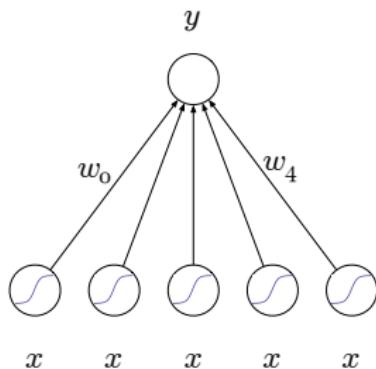
Every continuous function in the function space can be represented as a **linear combination*** of basis functions!

*Nothing new if you know Taylor expansion, Fourier transform, ...

Towards nonlinear systems

How do we find the optimal basis function ϕ ?

The above system could be graphically represented like this

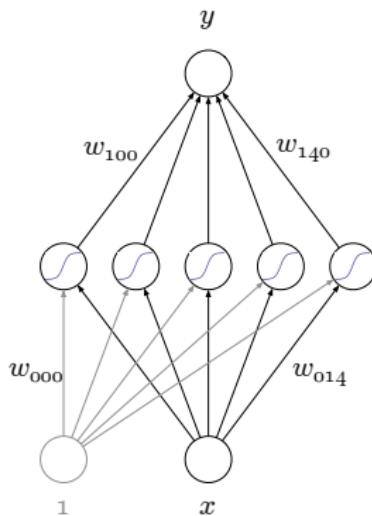


where the arrows represent the weights w and the circles the basis functions.

Why don't we let the system find the optimal ϕ ?

The multi-layered perceptron = neural network

Idea: we can extend the system by an additional layer, so that each basis function is itself a nonlinear function of a linear combination of the inputs:



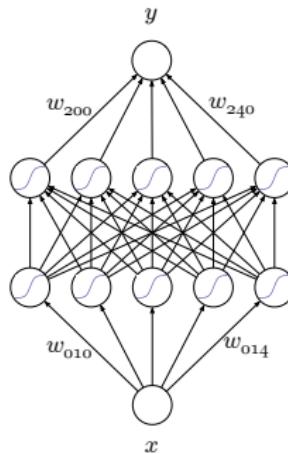
(for simplicity, the constant "1" is usually and from now on not depicted. But you always need it!)

We have generalised to $y(\mathbf{x}, \mathbf{w}_0, \mathbf{w}_1) = \mathbf{w}_1^T \mathbf{h}(\mathbf{w}_0^T \mathbf{x})$.

Note: in the neural network model, circles represent differentiable, nonlinear activation functions \mathbf{h} (*ReLU, sigmoid, tanh, ...*).

The deep neural network

We can continue adding more hidden layers



and get a deep neural network: $y(\mathbf{x}, \mathbf{w}) = \mathbf{w}_2^T h(\mathbf{w}_1^T h(\mathbf{w}_0^T \mathbf{x}))$.

This reduces the problem to finding the optimal \mathbf{w} !

How do we find \mathbf{w} ?

Usually we minimise the mean squared error (MSE):

$$\mathcal{L}(\mathbf{w}) = \sum_{n=1}^N (z_n - y(x_n, \mathbf{w}))^2$$

There is one difference w.r.t. linear regression: $\mathcal{L}(\mathbf{w})$ is no longer convex!

How can this be minimised? The minimum is located where its gradient is 0. So one typically minimises by using the gradient:

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} - \alpha \nabla \mathcal{L}(\mathbf{w})$$

How do we compute the gradient $\nabla \mathcal{L}$? Backpropagation does this by using the chain rule:

$$f(x) = u(v(x)) \quad \rightarrow \quad f'(x) = u'(v(x)) \cdot v'(x).$$

Algorithm for backprop (“on-line” aka “stochastic” learning)

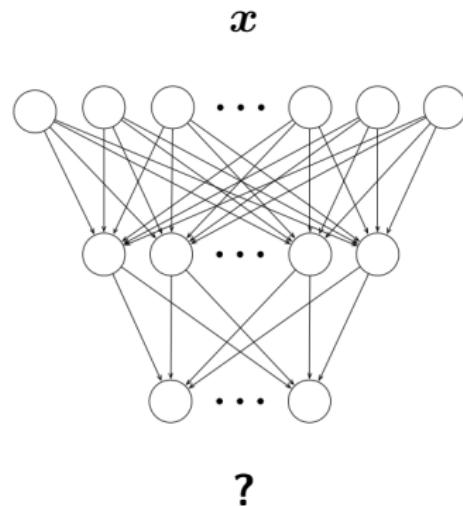
backpropagation algorithm:

```
initialise the weights
repeat
for each training sample  $(x, z)$  do
begin
compute  $\mathbf{o} = y(\mathbf{w}, \mathbf{x})$  (forward pass)
calculate residual  $\delta_{kj} = z - o$  at the output units
for all k:
propagate  $\delta_{kj}$  back one layer by  $\delta_{k-1,i} = \sum_j \delta_{kj} w_{k-1,i,j}$ 
update the weights using  $\partial \mathcal{L} / \partial w_{kij} = \delta_{kj} \phi'(\cdot) x_i$ 
end
(this is called one epoch)
until stopping criterion satisfied
```

Unsupervised learning with neural networks

Problem: often we have unlabelled data

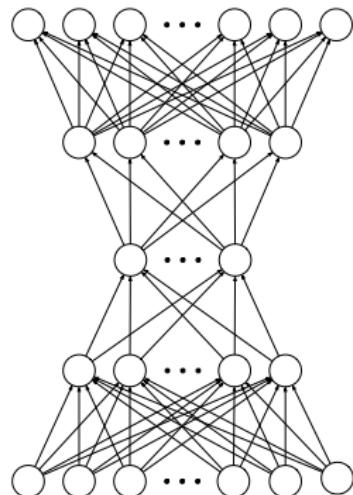
1. we have many data $x \dots$
2. but no the related output z .



Autoencoder networks

Idea: find compact representation of inputs (unsupervised!) by

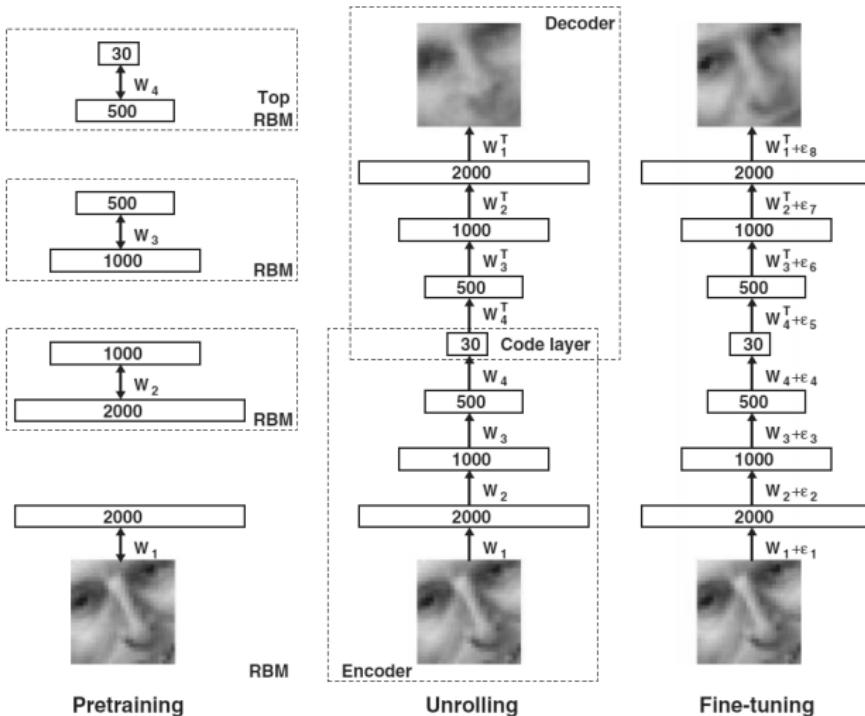
1. using a second neural network as decoder;
2. letting the neural network compute $y(x) = x$.



- ▶ $\mathcal{L}_{\text{AE}}(w) = \|x - y(x, w)\|^2$
- ▶ middle layer (“latents”) usually has fewer neurons
- ▶ latent representation $z = \text{compact code for } x$

These networks make a compact representation of data (*dimensionality reduction*).

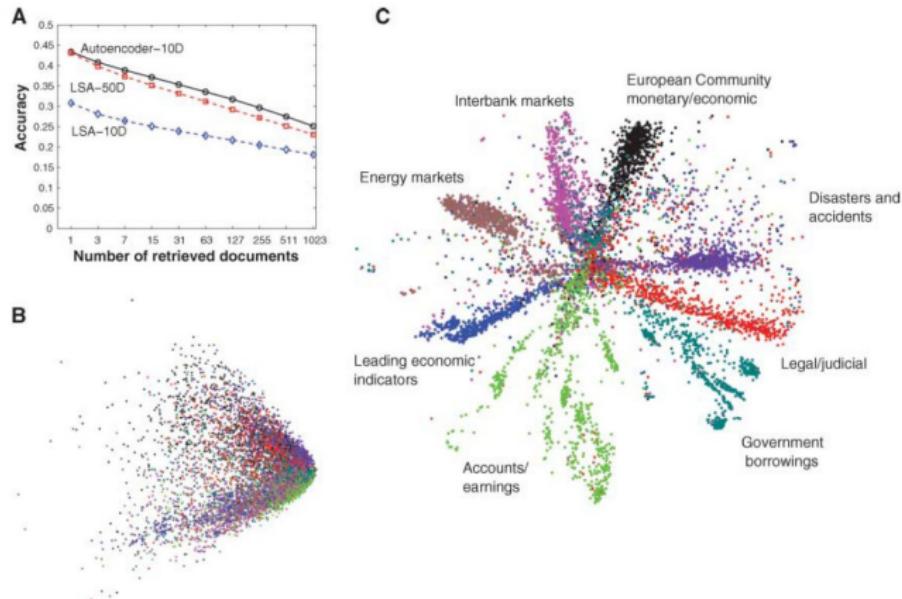
Autoencoder example: face compression



Hinton, G. E. and Salakhutdinov, R. R. Reducing the Dimensionality of Data with Neural Networks. *Science*, 2006.

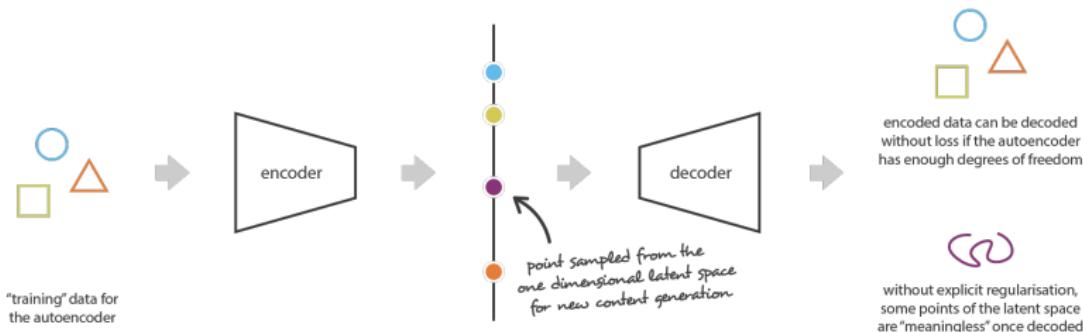
Autoencoder example: text categorisation

Fig. 4. (A) The fraction of retrieved documents in the same class as the query when a query document from the test set is used to retrieve other test set documents, averaged over all 402,207 possible queries. (B) The codes produced by two-dimensional LSA. (C) The codes produced by a 2000-500-250-125-2 autoencoder.



Limitations of autoencoders

An autoencoder cannot *generalise* from training data such that it is capable of generating new data!



The reason is we don't learn the underlying distribution of our data $p(x)$.

The probabilistic perspective is crucial

Why is the underlying distribution of our data interesting?

- ▶ Knowing $p(\mathbf{x})$ enables the generation of realistic objects, movements, etc.
- ▶ If our model learns, e.g., the syntax of a language, it can generate meaningful text.
- ▶ In case of sequential data, learning the physics/dynamics of a system allows the model to predict the future.

Generating photorealistic faces



Karras, T. et al. A Style-Based Generator Architecture for Generative Adversarial Networks. *CVPR*, 2019.

Generating coherent paragraphs of text

SYSTEM PROMPT
(HUMAN-WRITTEN)

Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today.

MODEL
COMPLETION
(MACHINE-
WRITTEN, SECOND
TRY)

The 19-year-old singer was caught on camera being escorted out of the store by security guards.

The singer was wearing a black hoodie with the label ‘Blurred Lines’ on the front and ‘Fashion Police’ on the back.

Scroll down for video

Shoplifting: Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today (pictured)

The singer was wearing a black hoodie with the label ‘Blurred Lines’ on the front and ‘Fashion Police’ on the back

The singer was also wearing a pair of black-rimmed glasses, a black jacket, black jeans and black sandals.

She was carrying a pair of black and white striped gloves and a small black bag.

Beyond visualisation: taking advantage of the latent representation

- ▶ Latent variables z are often interpreted as the *quintessential information*.
- ▶ What if we could utilise z to learn the underlying distribution of our data $p(\mathbf{x})$?
- ▶ Let's use the definition of latent-variable models as starting point:

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z} = \mathbb{E}_{p(\mathbf{z})}[p(\mathbf{x}|\mathbf{z})].$$

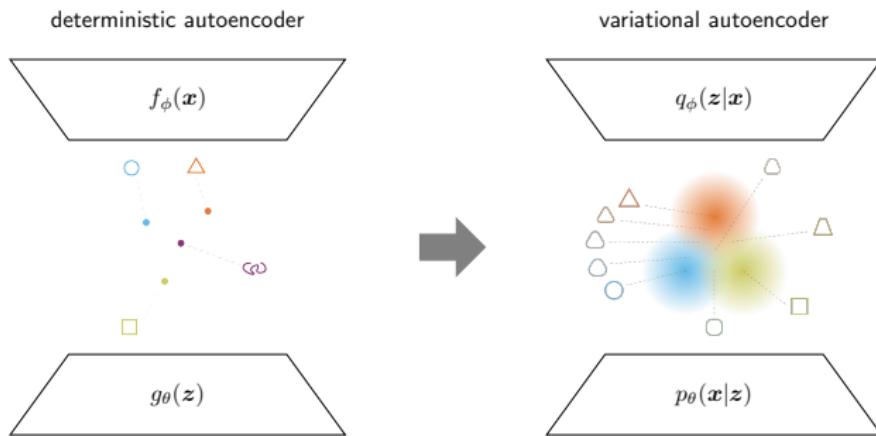
How do we find/learn $\mathbb{E}_{p(\mathbf{z})}[p(\mathbf{x}|\mathbf{z})]?$ *

*The prior $p(\mathbf{z})$ is usually defined in advance (e.g. as standard normal distribution).

Variational autoencoder—getting an intuition

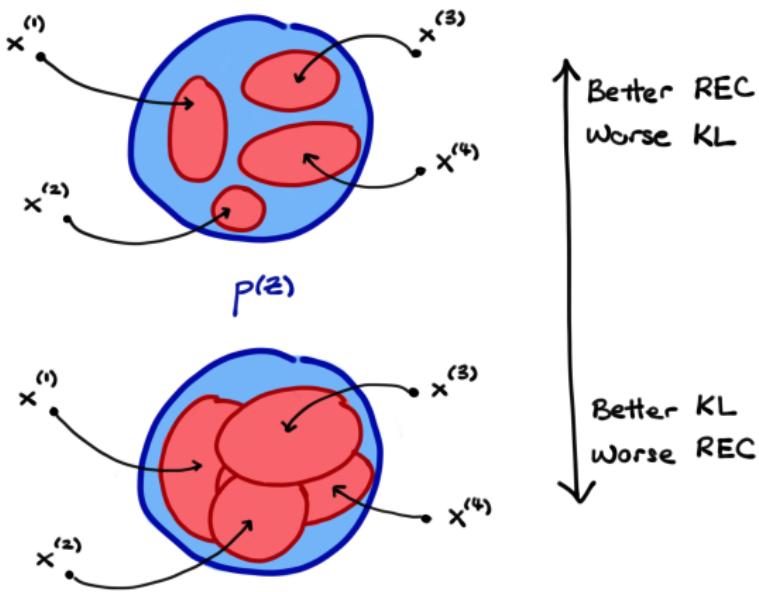
To learn $\mathbb{E}_{p(\mathbf{z})}[p(\mathbf{x}|\mathbf{z})]$, we introduce a variational distribution $q(\mathbf{z}|\mathbf{x})$.
The resulting objective function is known as the evidence lower bound:

$$\log p(\mathbf{x}) \geq \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{reconstruction}} - \underbrace{\text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))}_{\text{regularisation}} \equiv \text{ELBO}(\phi, \theta)$$



Note: ϕ and θ denote the parameters of the functions/distributions learned by the neural networks.

KL divergence as regulariser



Deriving the evidence lower bound

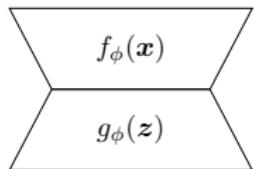
$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \\&= \log \int \frac{p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\&= \log \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\frac{p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \\&\stackrel{1}{\geq} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \\&= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})} \right] \\&= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})) \\&\quad \left(\stackrel{2}{\propto} - \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\underbrace{\frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}_{\theta}(\mathbf{z})\|^2}_{\text{reconstruction error}} \right] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})) \right)\end{aligned}$$

¹Jensen's inequality

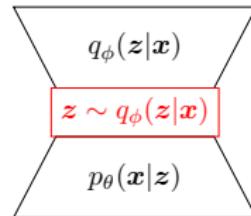
²If $p_{\theta}(\mathbf{x}|\mathbf{z})$ is a diagonal Gaussian with $\Sigma = \mathbb{I}$.

Deterministic vs. variational autoencoder—Gaussian case*

deterministic autoencoder



variational autoencoder



$$\nabla_{\phi, \theta} \mathcal{L}_{\text{AE}} = \nabla_{\phi, \theta} \| \mathbf{x} - g_\theta(f_\phi(\mathbf{x})) \|^2$$

$$\nabla_{\phi, \theta} \mathcal{L}_{\text{VAE}} = -\nabla_{\phi, \theta} \text{ELBO}$$

$$= \nabla_{\phi, \theta} \left[\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} \left[\frac{1}{2} \| \mathbf{x} - \boldsymbol{\mu}_\theta(\mathbf{z}) \|^2 \right] + \text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})) \right]$$

The ELBO rewards good reconstruction (first term), and comes with an inbuilt regulariser (the KL) that prevents collapsing to the deterministic autoencoder.

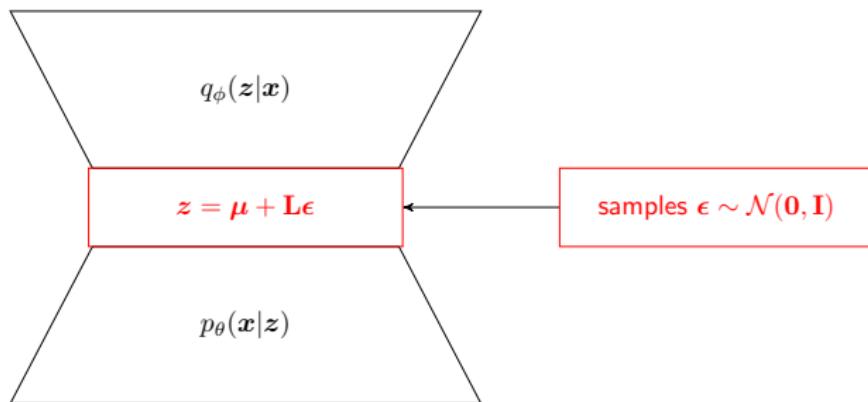
* In the VAE setting, $q_\phi(\mathbf{z}|\mathbf{x})$ and $p_\theta(\mathbf{x}|\mathbf{z})$ are typically specified to be diagonal Gaussian distributions.

One important detail: backpropagation through random sampling?

Solution: reparametrisation

$$\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \rightsquigarrow \quad \mathbf{z} = \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$$

This allows taking partial derivatives w. r. t. $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$.



The popular MNIST dataset

Consider the following data:

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

MNIST handwritten digits

<http://deeplearning.net/data/mnist/>

What might be the underlying distribution $p(\mathbf{x})$?

Our standard set of distributions is not expressive enough:



Learning process of a variational autoencoder

If this figure is not animated, you might want to try a different pdf reader.

For visualisation, the model only uses a two-dimensional latent space. A higher dimension would allow the model to achieve better results.

Some results from our lab: human motion

If this figure is not animated, you might want to try a different pdf reader.

Some results from our lab: pendulum (sequential data)

If this figure is not animated, you might want to try a different pdf reader.

Karl, M. et al. Deep Variational Bayes Filters: Unsupervised Learning of State Space Models from Raw Data. *ICLR*, 2017.

What we learned

- ▶ Neural networks are function approximators.
- ▶ Two neural networks can be combined to an autoencoder for learning a compact/latent representation of our data in an unsupervised fashion.
- ▶ The variational autoencoder—a probabilistic method—can learn the underlying distribution of our data by means of latent variables.