

Lecture 11. Autoencoders, Deep generative models

COMP90051 Statistical Machine Learning

Semester 1, 2021
Trevor Cohn



THE UNIVERSITY OF
MELBOURNE

This lecture

- Autoencoders
 - * Learning efficient coding
 - * Methods for autoencoding: sparse, denoising, contractive
 - * Use in pre-training pipelines
- Deep generative models
 - * Variational autoencoder (VAE)

Autoencoder

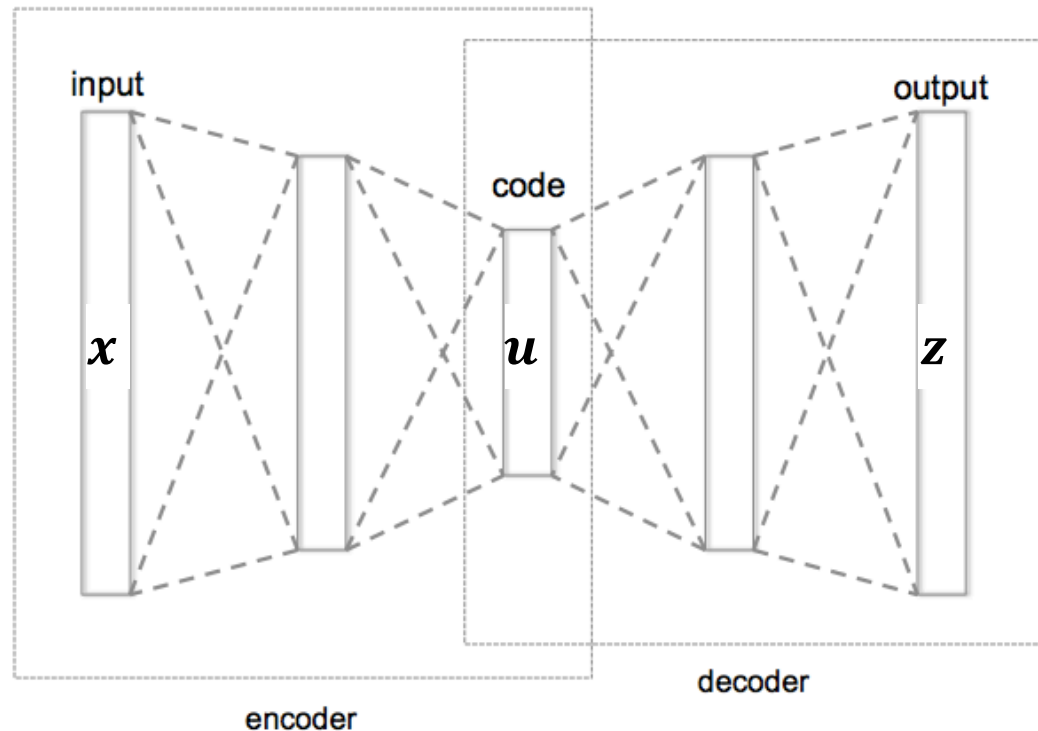
An ANN training setup that can be used
for unsupervised learning, initialisation,
or just efficient coding

Autoencoding idea

- Supervised learning:
 - * Univariate regression: predict y from x
 - * Multivariate regression: predict \mathbf{y} from \mathbf{x}
- Unsupervised learning: explore data $\mathbf{x}_1, \dots, \mathbf{x}_n$
 - * No response variable
- For each \mathbf{x}_i set $\mathbf{y}_i \equiv \mathbf{x}_i$
- Train an ANN to predict \mathbf{y}_i from \mathbf{x}_i i.e., model $p(\mathbf{x}|\mathbf{x})$
- Pointless?

Autoencoder topology

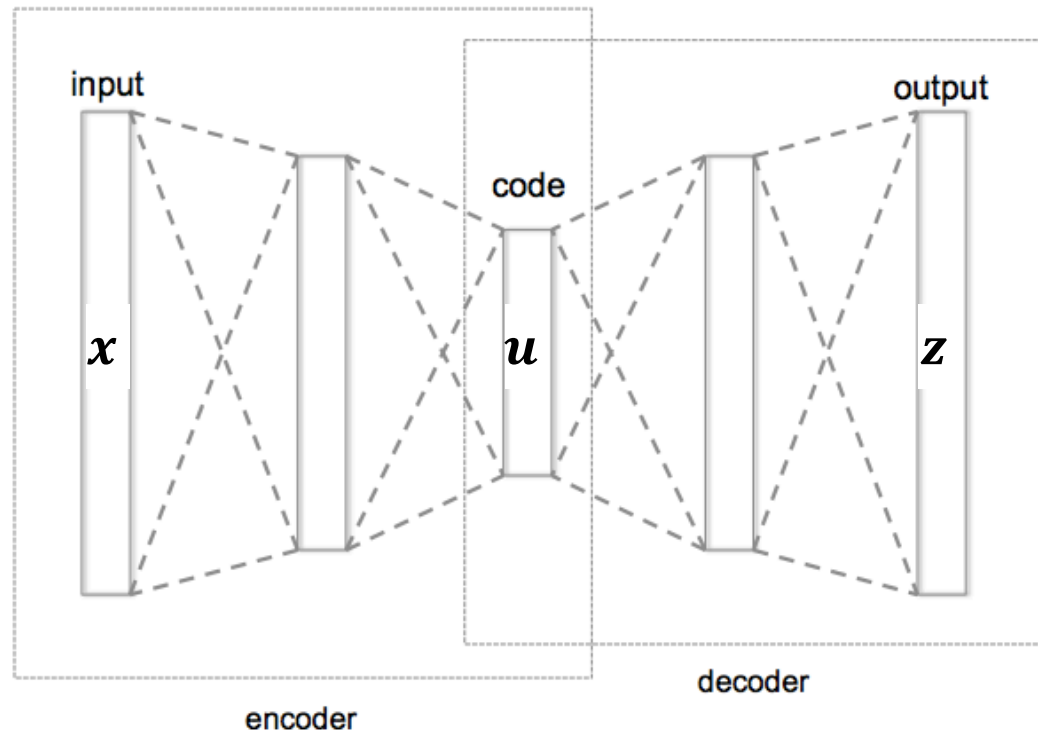
- Given data without labels $\mathbf{x}_1, \dots, \mathbf{x}_n$, set $\mathbf{y}_i \equiv \mathbf{x}_i$ and train an ANN to predict $\mathbf{z}(\mathbf{x}_i) \approx \mathbf{x}_i$
- Set **bottleneck** layer \mathbf{u} in middle “thinner” than input



adapted from: Chervinskii at
Wikimedia Commons (CC4)

Introducing the bottleneck

- Suppose you managed to train a network that gives a good **restoration** of the original signal $\mathbf{z}(\mathbf{x}_i) \approx \mathbf{x}_i$
- This means that the data structure can be effectively described (**encoded**) by a lower dimensional representation \mathbf{u}



adapted from: Chervinskii at
Wikimedia Commons (CC4)

Under-/Over-completeness

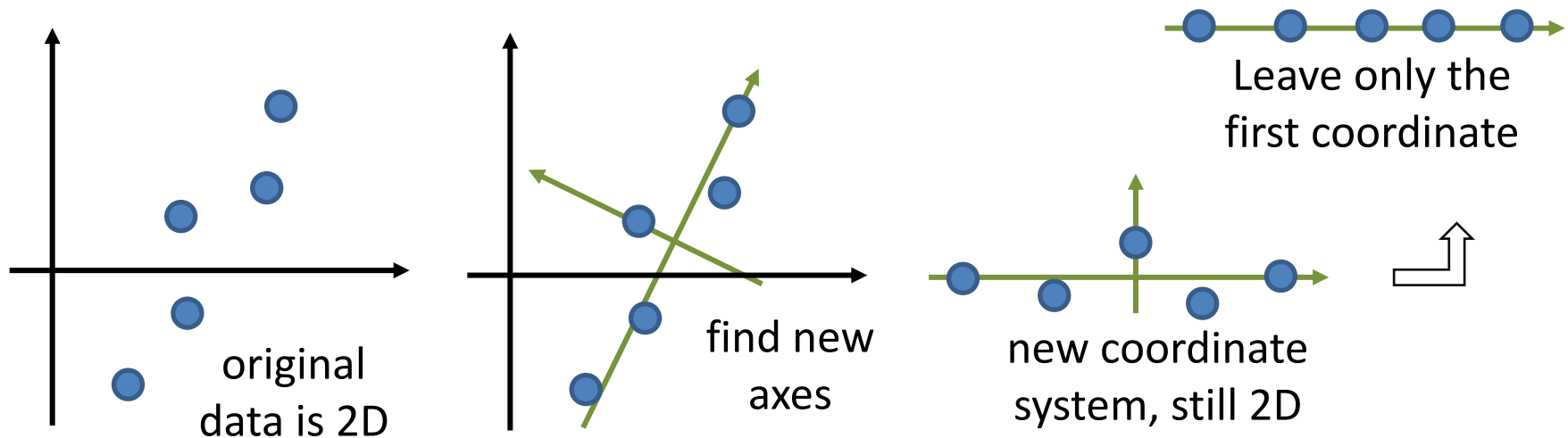
- Manner of bottleneck gives rise to:
 - * **undercomplete**: model with thinner bottleneck than input forced to generalise
 - * **overcomplete**: wider bottleneck than input, can just “copy” input directly to output
- Even undercomplete models can learn trivial codes, given complex non-linear encoder and decoder
- Various methods to ensure learning

Dimensionality reduction

- Autoencoders can be used for **compression** and **dimensionality reduction** via a non-linear transformation
- Closely related to principal component analysis (PCA)...

Principal component analysis

- Principal component analysis (PCA) is a popular method for dimensionality reduction and data analysis in general
- Given a dataset $x_1, \dots, x_n, x_i \in \mathbf{R}^m$, PCA aims to find a new coordinate system such that most of the **variance is concentrated** along the first coordinate, then most of the remaining variance along the second coordinate, etc.
- Dimensionality reduction is then based on **discarding coordinates** except the first $l < m$



PCA: Solving the optimisation

- PCA aims to find \mathbf{p}_1 that maximises $\mathbf{p}_1' \mathbf{\Sigma}_X \mathbf{p}_1$, subject to $\|\mathbf{p}_1\| = \mathbf{p}_1' \mathbf{p}_1 = 1$
- Constrained \rightarrow Lagrange multipliers. Introduce multiplier λ_1 ; set derivatives of Lagrangian to zero, solve
- $L = \mathbf{p}_1' \mathbf{\Sigma}_X \mathbf{p}_1 - \lambda_1 (\mathbf{p}_1' \mathbf{p}_1 - 1)$
- $\frac{\partial L}{\partial \mathbf{p}_1} = 2\mathbf{\Sigma}_X \mathbf{p}_1 - 2\lambda_1 \mathbf{p}_1 = 0$
- $\mathbf{\Sigma}_X \mathbf{p}_1 = \lambda_1 \mathbf{p}_1$
- Precisely defines \mathbf{p}_1 as an **eigenvector** with λ_1 being the corresponding **eigenvalue**

PCA vs Autoencoding

- If you use linear activation functions and only one hidden layer, then the setup becomes almost that of **Principal Component Analysis (PCA)**
 - * PCA finds orthonormal basis where axes are aligned to capture maximum data variation
 - * ANN might find a different solution, doesn't use eigenvalues (directly)

Mini-summary

- Autoencoding as modelling $p(\mathbf{x}|\mathbf{x})$
- Notion of under- and over-completeness
- Relationship to PCA

Autoencoding techniques

Methods for learning Autoencoders

- Sparse autoencoders
 - * includes regularisation over the code layer
- Denoising autoencoders
 - * learn to recover true data given noise-corrupted input
- Contractive autoencoders
 - * explicit penalty term to ensure robustness to local changes in input
- In all cases, manage the tension between compact code and accurate reconstruction

Sparse Autoencoders

- Apply constraint to the hidden code
 - * Use regularisation penalty in training objective, e.g.,

$$\lambda \sum_i |u_i|$$

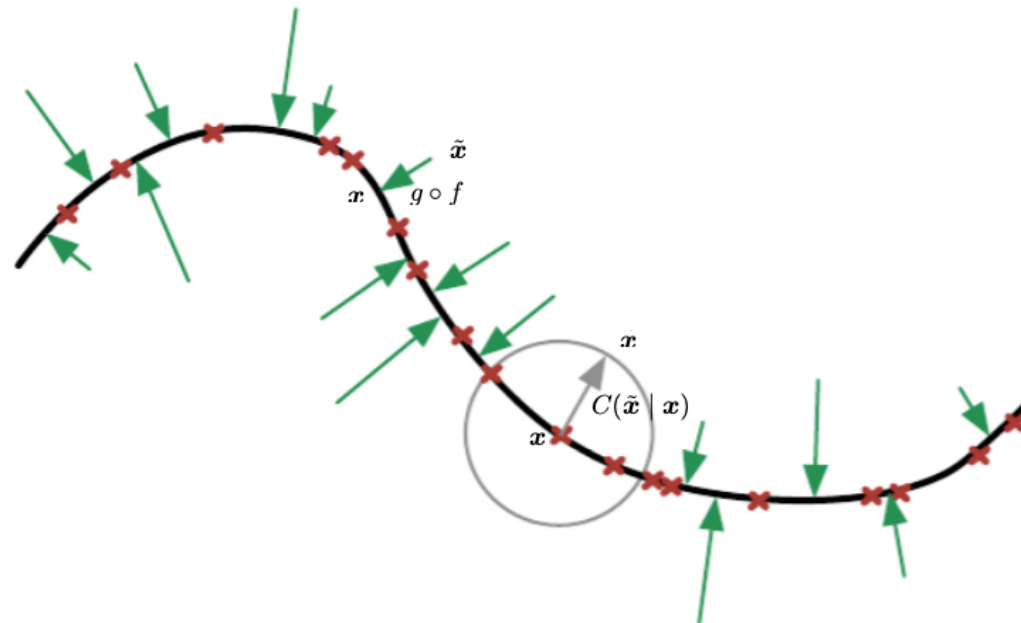
- * May need to use with over-complete hidden layer, and multi-layer encoder/decoder
 - * Sparsity can simplify data analysis
- Can be thought of as a kind of as prior

Denoising Autoencoder

- Corrupt input data to the autoencoder, but applying noise (masking, additive whitenoise, etc)
 - * Now modelling $p(\mathbf{x}|\tilde{\mathbf{x}})$ where $\tilde{\mathbf{x}}$ is noised input
 - * The model must recover corrupted parts of the input, thus must learn underlying structure of the data
- Popular “BERT” model for NLP uses this idea to learn word and sentence representations
 - * Based on masking of words in a transformer neural network over strings

AE and Manifold Learning

- Most of the input space does not correspond to input data, data lives on lower dim “manifold”
 - * DAE learns training data manifold, by mapping nearby points onto the manifold

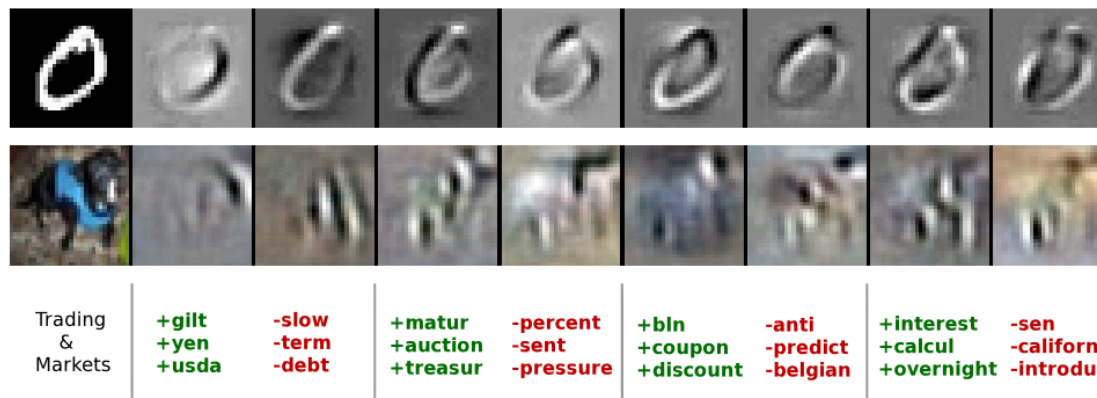


Contractive Autoencoder

- Include penalty term over derivatives wrt input

$$\lambda \sum_i ||\nabla_x u_i||^2$$

- * local changes to input \mathbf{x} have limited effect on code \mathbf{u}
- * tangent vectors to the manifold encode data variation



Rifai, S., Dauphin, Y., Vincent, P., Bengio, Y., and Muller, X. (2011c). The manifold tangent classifier. In NIPS'2011.

Uses of Autoencoders

- Data visualisation & clustering
 - * Unsupervised first step towards understanding properties of the data
- As a feature representation
 - * Allowing the use of off-the-shelf ML methods, applied to much smaller and informative representations of input
- Pre-training of deep models
 - * Warm-starting training by initialising model weights with encoder parameters
 - * In some fields like vision, mostly replaced with supervised pre-training on very large datasets

Mini-summary

- Autoencoding methods
 - * Sparse
 - * Denoising
 - * Contractive
- Relation to manifold learning
- Uses of autoencoders

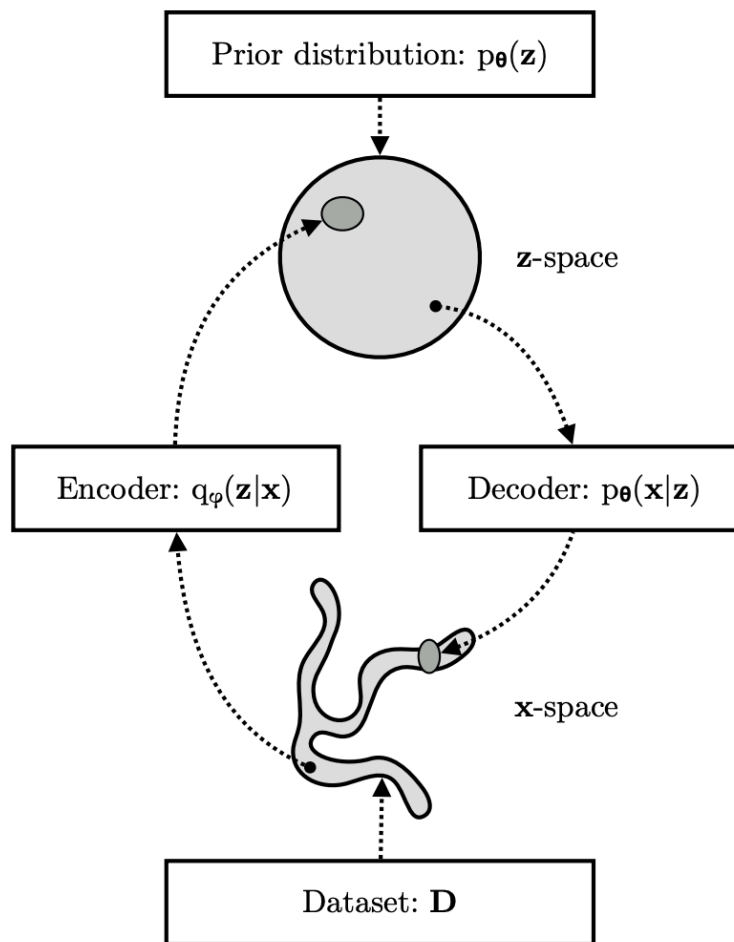
Variational “autoencoder”

A form of deep generative model
whose inference method resembles
that of an autoencoder

Deep generative models

- Often need a model which can *generate inputs*
 - * trained to model $p(\mathbf{x})$ directly. I.e., unsupervised model trained as density estimator
 - * sometimes we also model a label, $p(\mathbf{x}, y)$ — we will return to this in a later lecture
- Use a probabilistic model with a **latent random variable \mathbf{z}** to explain data variation
 - * $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$ (*sum denotes integral or summation*)
 - * latent variable to capture structure in data, such as clusters, style, etc
 - * framework includes several models, e.g., RBMs, DBMs etc

Structure of Variational Autoencoder



Variational Autoencoder

- Given generative model, $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$, how to fit to data?

* $\log p(\mathbf{x})$

$$= \log \sum_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$$

$$= \log \sum_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \frac{q(\mathbf{z}|\mathbf{x})}{q(\mathbf{z}|\mathbf{x})}$$

$$= \log E_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \left[p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \frac{1}{q(\mathbf{z}|\mathbf{x})} \right]$$

$$\geq E_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \left[\log p(\mathbf{x}|\mathbf{z}) \frac{p(\mathbf{z})}{q(\mathbf{z}|\mathbf{x})} \right]$$

assume model factorises

introduce q distribution
(chosen carefully)

use Jensen's inequality
to move log inside
expectation

- * This is known as the “evidence lower bound” or ELBO. We seek to maximise this quantity.

VAE (cont.)

- Have to choose distributions, most commonly:
 - * prior $p(\mathbf{z})$ a unit Gaussian, $N(0,1)$
 - * $q(\mathbf{z}|\mathbf{x})$ also a Gaussian, $N(\mu(\mathbf{x}), \sigma^2(\mathbf{x}))$ as neural network with vector outputs, and parameters ϕ
 - * likelihood $p(\mathbf{x}|\mathbf{z})$ another neural network, parameters θ
- Optimisation problem is now

$$\arg \max_{\theta, \phi} E_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \left[\log p(\mathbf{x}|\mathbf{z}) + \log \frac{p(\mathbf{z})}{q(\mathbf{z}|\mathbf{x})} \right]$$

reconstruction term

Kullback-Leibler divergence
 $-D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$

Optimising the VAE

- Intractable objective due to expectation
 - * approximate expectation with single sample $\hat{\mathbf{z}} \sim q(\mathbf{z}|\mathbf{x})$
 - * but sampling is not a differentiable operator, can't backpropagate gradients (to learn ϕ)
- Reparameterisation trick
 - * sample instead $\epsilon \sim N(0,1)$, then rescale and offset,
$$\hat{\mathbf{z}} = \mu(\mathbf{x}) + \hat{\epsilon} \odot \sigma^2(\mathbf{x})$$
 - * back-propagation now side-steps sampling
 - * N.b., trick only works for a limited set of distributions

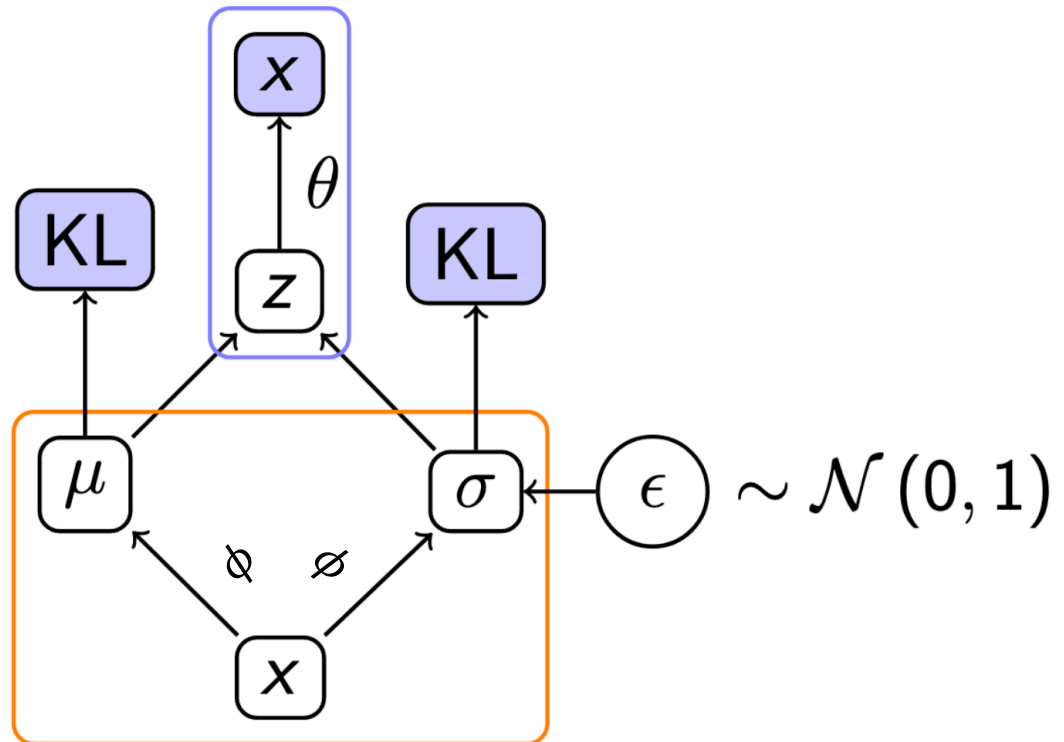
VAE as a Computation Graph

generation model

a.k.a., “decoder”

inference model

a.k.a., “encoder”

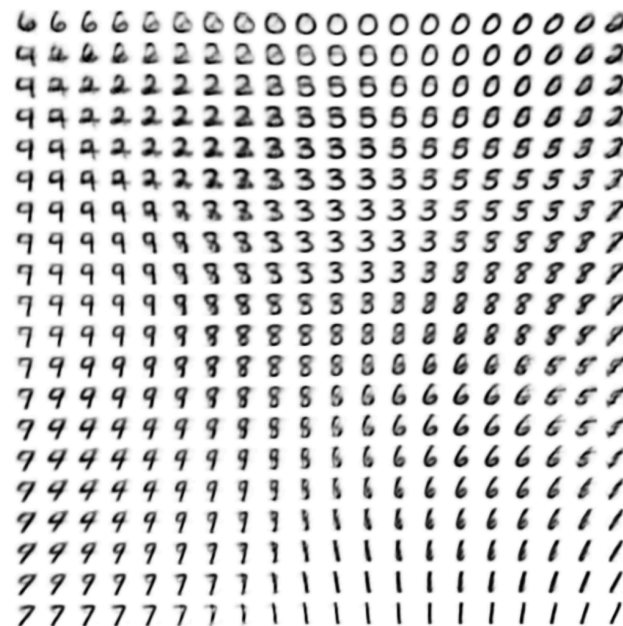


Generating from a VAE

- Simply sample from prior, $\hat{\mathbf{z}} \sim p(\mathbf{z})$, then sample from decoder, $\hat{\mathbf{x}} \sim p(\mathbf{x}|\hat{\mathbf{z}})$
- Examples



(a) Learned Frey Face manifold



(b) Learned MNIST manifold

Kingma, Diederik P., and Max Welling. "Auto-encoding variational bayes." *ICLR 2014*

Figure 4: Visualisations of learned data manifold for generative models with two-dimensional latent space, learned with AEVB. Since the prior of the latent space is Gaussian, linearly spaced coordinates on the unit square were transformed through the inverse CDF of the Gaussian to produce values of the latent variables \mathbf{z} . For each of these values \mathbf{z} , we plotted the corresponding generative $p_{\theta}(\mathbf{x}|\mathbf{z})$ with the learned parameters θ .

This lecture

- Autoencoders
 - * Learning efficient coding
 - * Methods for autoencoding
 - * Their uses
- Variational autoencoders
- **Workshops Week #6**: Convolution & Autoencoding
- Next lectures: Recurrent neural networks