

# Probabilistic modelling for NLP powered by deep learning

Wilker Aziz  
University of Amsterdam

April 3, 2018

# Outline

- 1 Deep generative models
- 2 Variational inference
- 3 Variational auto-encoder

# Problems

## Supervised problems

*“learn a distribution over **observed** data”*

- sentences in natural language
- images, ...

# Problems

## Supervised problems

*“learn a distribution over **observed** data”*

- sentences in natural language
- images, ...

## Unsupervised problems

*“learn a distribution over **observed** and **unobserved** data”*

- sentences in natural language + parse trees
- images + bounding boxes, ...

# Supervised problems

We have data  $x^{(1)}, \dots, x^{(N)}$  e.g.

- sentences, images, ...

generated by some **unknown** procedure

# Supervised problems

We have data  $x^{(1)}, \dots, x^{(N)}$  e.g.

- sentences, images, ...

generated by some **unknown** procedure  
which we assume can be captured by a probabilistic model

# Supervised problems

We have data  $x^{(1)}, \dots, x^{(N)}$  e.g.

- sentences, images, ...

generated by some **unknown** procedure

which we assume can be captured by a probabilistic model

- with **known** probability (mass/density) function e.g.

$$\underbrace{X \sim \text{Cat}(\pi_1, \dots, \pi_K)}_{\text{e.g. nationality}}$$

or

$$\underbrace{X \sim \mathcal{N}(\mu, \sigma^2)}_{\text{e.g. height}}$$

# Supervised problems

We have data  $x^{(1)}, \dots, x^{(N)}$  e.g.

- sentences, images, ...

generated by some **unknown** procedure

which we assume can be captured by a probabilistic model

- with **known** probability (mass/density) function e.g.

$$\underbrace{X \sim \text{Cat}(\pi_1, \dots, \pi_K)}_{\text{e.g. nationality}}$$

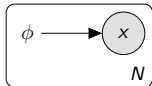
or

$$\underbrace{X \sim \mathcal{N}(\mu, \sigma^2)}_{\text{e.g. height}}$$

**estimate parameters** that assign maximum likelihood to observations



# Multiple problems, same language



(Conditional) Density estimation

	Side information ( $\phi$ )	Observation ( $x$ )
Parsing	a sentence	parse tree
Translation	a sentence in French	translation in English
Captioning	an image	caption in English
Entailment	a text and hypothesis	entailment relation

# Where does deep learning kick in?

Let  $\phi$  be all side information available  
e.g. deterministic *inputs/features*

# Where does deep learning kick in?

Let  $\phi$  be all side information available  
e.g. deterministic *inputs/features*

Have neural networks predict parameters of our probabilistic model

$$X|\phi \sim \text{Cat}(\pi_{\mathbf{w}}(\phi)) \quad \text{or} \quad X|\phi \sim \mathcal{N}(\mu_{\mathbf{w}}(\phi), \sigma_{\mathbf{w}}(\phi)^2)$$

# Where does deep learning kick in?

Let  $\phi$  be all side information available  
e.g. deterministic *inputs/features*

Have neural networks predict parameters of our probabilistic model

$$X|\phi \sim \text{Cat}(\pi_{\mathbf{w}}(\phi)) \quad \text{or} \quad X|\phi \sim \mathcal{N}(\mu_{\mathbf{w}}(\phi), \sigma_{\mathbf{w}}(\phi)^2)$$

and proceed to **estimate parameters**  $\mathbf{w}$  of the NNs

# NN as efficient parametrisation

From the statistical point of view NNs do not generate data

- they parametrise distributions that *by assumption* govern data
- compact and efficient way to map from complex side information to parameter space

# NN as efficient parametrisation

From the statistical point of view NNs do not generate data

- they parametrise distributions that *by assumption* govern data
- compact and efficient way to map from complex side information to parameter space

Prediction is done by a decision rule outside the statistical model

- e.g. beam search

# MLE via gradient-based optimisation

The probability of an observation  $X = x$  is given by some **differentiable** probability function  $p(x)$

- the parameters of which are predicted by  $f_w$   
(*also differentiable*)

# MLE via gradient-based optimisation

The probability of an observation  $X = x$  is given by some **differentiable** probability function  $p(x)$

- the parameters of which are predicted by  $f_w$   
*(also differentiable)*

Example:  $K$  classes

$$p(x) = \text{Cat}(X = x | \underbrace{\pi_1^K := f_w(\phi)}_{\text{class probabilities}}) = \prod_{i=1}^K \pi_i^{[x=i]}$$



# MLE via gradient-based optimisation

The probability of an observation  $X = x$  is given by some **differentiable** probability function  $p(x)$

- the parameters of which are predicted by  $f_w$   
*(also differentiable)*

Example:  $K$  classes

$$p(x) = \text{Cat}(X = x | \underbrace{\pi_1^K := f_w(\phi)}_{\text{class probabilities}}) = \prod_{i=1}^K \pi_i^{[x=i]}$$

Given a dataset of i.i.d. observations, SGD gives us a local optimum of the log-likelihood

# DL in NLP recipe

## Maximum likelihood estimation

- tells you which **loss** to optimise  
(i.e. negative log-likelihood)

## DL in NLP recipe

Maximum likelihood estimation

- tells you which **loss** to optimise  
(i.e. negative log-likelihood)

Automatic differentiation (*backprop*)

- chain rule of derivatives: “give me a tractable forward pass  
and I will give you **gradients**”

## DL in NLP recipe

### Maximum likelihood estimation

- tells you which **loss** to optimise  
(i.e. negative log-likelihood)

### Automatic differentiation (*backprop*)

- chain rule of derivatives: “give me a tractable forward pass and I will give you **gradients**”

### Stochastic optimisation powered by backprop

- general purpose gradient-based optimisers

# Tractability is central

- Likelihood gives us a differentiable objective to optimise for
- but we need to stick with **tractable** likelihood functions

# When do we have intractable likelihood?

## Unsupervised problems

assessing the likelihood requires marginalisation of latent variables

# When do we have intractable likelihood?

## Unsupervised problems

assessing the likelihood requires **marginalisation of latent variables**

- too many forward passes

$$p(x) = \sum_{c=1}^K \text{Cat}(c|\pi_1, \dots, \pi_K) \underbrace{\mathcal{N}(x|\mu_w(c), \sigma_w(c)^2)}_{\text{forward pass}}$$

# When do we have intractable likelihood?

## Unsupervised problems

assessing the likelihood requires **marginalisation of latent variables**

- too many forward passes

$$p(x) = \sum_{c=1}^K \text{Cat}(c|\pi_1, \dots, \pi_K) \underbrace{\mathcal{N}(x|\mu_w(c), \sigma_w(c)^2)}_{\text{forward pass}}$$

- even infinitely many

$$p(x) = \int \mathcal{N}(z|0, I) \underbrace{\text{Cat}(x|\pi_w(z))}_{\text{forward pass}} dz$$



# Deep generative models

Joint distribution with **deep observation model**

$$p_{\theta}(x, z) = \underbrace{p(z)}_{\text{prior}} \underbrace{p_{\theta}(x|z)}_{\text{likelihood}}$$

mapping from latent variable  $z$  to  $p(x|z)$  is a NN with parameters  $\theta$

# Deep generative models

Joint distribution with **deep observation model**

$$p_{\theta}(x, z) = \underbrace{p(z)}_{\text{prior}} \underbrace{p_{\theta}(x|z)}_{\text{likelihood}}$$

mapping from latent variable  $z$  to  $p(x|z)$  is a NN with parameters  $\theta$

Marginal likelihood (or evidence)

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

**intractable** in general

# Gradient

Exact gradient is intractable

$$\nabla_{\theta} \log p_{\theta}(x)$$

Exact gradient is intractable

$$\nabla_{\theta} \log p_{\theta}(x) = \nabla_{\theta} \log \underbrace{\int p_{\theta}(x, z) \, dz}_{\text{marginal}}$$

Exact gradient is intractable

$$\begin{aligned}\nabla_{\theta} \log p_{\theta}(x) &= \nabla_{\theta} \log \underbrace{\int p_{\theta}(x, z) \, dz}_{\text{marginal}} \\ &= \underbrace{\frac{1}{\int p_{\theta}(x, z) \, dz} \int \nabla_{\theta} p_{\theta}(x, z) \, dz}_{\text{chain rule}}\end{aligned}$$

Exact gradient is intractable

$$\begin{aligned}\nabla_{\theta} \log p_{\theta}(x) &= \nabla_{\theta} \log \underbrace{\int p_{\theta}(x, z) \, dz}_{\text{marginal}} \\&= \frac{1}{\underbrace{\int p_{\theta}(x, z) \, dz}_{\text{chain rule}}} \underbrace{\int \nabla_{\theta} p_{\theta}(x, z) \, dz}_{\text{chain rule}} \\&= \frac{1}{p_{\theta}(x)} \int \underbrace{p_{\theta}(x, z) \nabla_{\theta} \log p_{\theta}(x, z)}_{\text{log-identity for derivatives}} \, dz\end{aligned}$$

Exact gradient is intractable

$$\begin{aligned}\nabla_{\theta} \log p_{\theta}(x) &= \nabla_{\theta} \log \underbrace{\int p_{\theta}(x, z) \, dz}_{\text{marginal}} \\&= \underbrace{\frac{1}{\int p_{\theta}(x, z) \, dz} \int \nabla_{\theta} p_{\theta}(x, z) \, dz}_{\text{chain rule}} \\&= \frac{1}{p_{\theta}(x)} \int \underbrace{p_{\theta}(x, z) \nabla_{\theta} \log p_{\theta}(x, z)}_{\text{log-identity for derivatives}} \, dz \\&= \mathbb{E}_{p_{\theta}(z|x)} [\nabla_{\theta} \log p_{\theta}(x, Z)]\end{aligned}$$

Exact gradient is intractable

$$\begin{aligned}\nabla_{\theta} \log p_{\theta}(x) &= \nabla_{\theta} \log \underbrace{\int p_{\theta}(x, z) \, dz}_{\text{marginal}} \\&= \frac{1}{\underbrace{\int p_{\theta}(x, z) \, dz}_{\text{chain rule}}} \int \nabla_{\theta} p_{\theta}(x, z) \, dz \\&= \frac{1}{p_{\theta}(x)} \int \underbrace{p_{\theta}(x, z) \nabla_{\theta} \log p_{\theta}(x, z)}_{\text{log-identity for derivatives}} \, dz \\&= \mathbb{E}_{p_{\theta}(z|x)} [\nabla_{\theta} \log p_{\theta}(x, Z)]\end{aligned}$$

MC estimate of gradient requires sampling from posterior

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

unavailable due to the intractability of the marginal



# Summary

- We like probabilistic models because can make explicit modelling assumptions

# Summary

- We like probabilistic models because can make explicit modelling assumptions
- We want complex observation models parameterised by NNs

# Summary

- We like probabilistic models because can make explicit modelling assumptions
- We want complex observation models parameterised by NNs
- But we cannot use backprop for parameter estimation

# Summary

- We like probabilistic models because can make explicit modelling assumptions
- We want complex observation models parameterised by NNs
- But we cannot use backprop for parameter estimation

We need **approximate inference** techniques!

# Outline

- 1 Deep generative models
- 2 Variational inference
- 3 Variational auto-encoder

# The Basic Problem

The marginal likelihood

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

is generally **intractable**, which prevents us from computing quantities that depend on the posterior  $p(z|x)$

- e.g. gradients in MLE
- e.g. predictive distribution in Bayesian modelling

# Strategy

Accept that  $p(z|x)$  is not computable.

# Strategy

Accept that  $p(z|x)$  is not computable.

- approximate it by an auxiliary distribution  $q(z|x)$  that is computable
- choose  $q(z|x)$  as close as possible to  $p(z|x)$  to obtain a faithful approximation



# Evidence lowerbound

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

## Evidence lowerbound

$$\begin{aligned}\log p(x) &= \log \int p(x, z) dz \\ &= \log \int \textcolor{red}{q(z|x)} \frac{p(x, z)}{\textcolor{red}{q(z|x)}} dz\end{aligned}$$

## Evidence lowerbound

$$\begin{aligned}\log p(x) &= \log \int p(x, z) dz \\ &= \log \int q(z|x) \frac{p(x, z)}{q(z|x)} dz \\ &= \log \left( \mathbb{E}_{q(z|x)} \left[ \frac{p(x, z)}{q(z|x)} \right] \right)\end{aligned}$$

# Evidence lowerbound

$$\begin{aligned}
 \log p(x) &= \log \int p(x, z) dz \\
 &= \log \int q(z|x) \frac{p(x, z)}{q(z|x)} dz \\
 &= \log \left( \mathbb{E}_{q(z|x)} \left[ \frac{p(x, z)}{q(z|x)} \right] \right) \\
 &\geq \underbrace{\mathbb{E}_{q(z|x)} \left[ \log \frac{p(x, z)}{q(z|x)} \right]}_{\text{ELBO}}
 \end{aligned}$$

# Evidence lowerbound

$$\begin{aligned}
 \log p(x) &= \log \int p(x, z) dz \\
 &= \log \int q(z|x) \frac{p(x, z)}{q(z|x)} dz \\
 &= \log \left( \mathbb{E}_{q(z|x)} \left[ \frac{p(x, z)}{q(z|x)} \right] \right) \\
 &\geq \underbrace{\mathbb{E}_{q(z|x)} \left[ \log \frac{p(x, z)}{q(z|x)} \right]}_{\text{ELBO}} \\
 &= \mathbb{E}_{q(z|x)} [\log p(x, z)] + \mathbb{H}(q(z|x))
 \end{aligned}$$

# An approximate posterior

$$\log p(x) \geq \underbrace{\mathbb{E}_{q(z|x)} \left[ \log \frac{p(x, z)}{q(z|x)} \right]}_{\text{ELBO}}$$

# An approximate posterior

$$\begin{aligned}\log p(x) &\geq \underbrace{\mathbb{E}_{q(z|x)} \left[ \log \frac{p(x, z)}{q(z|x)} \right]}_{\text{ELBO}} \\ &= \mathbb{E}_{q(z|x)} \left[ \log \frac{p(z|x)p(x)}{q(z|x)} \right]\end{aligned}$$

# An approximate posterior

$$\begin{aligned}
 \log p(x) &\geq \underbrace{\mathbb{E}_{q(z|x)} \left[ \log \frac{p(x, z)}{q(z|x)} \right]}_{\text{ELBO}} \\
 &= \mathbb{E}_{q(z|x)} \left[ \log \frac{p(z|x)p(x)}{q(z|x)} \right] \\
 &= \int q(z|x) \log \frac{p(z|x)}{q(z|x)} dz + \underbrace{\log p(x)}_{\text{constant}}
 \end{aligned}$$



# An approximate posterior

$$\begin{aligned}
 \log p(x) &\geq \underbrace{\mathbb{E}_{q(z|x)} \left[ \log \frac{p(x, z)}{q(z|x)} \right]}_{\text{ELBO}} \\
 &= \mathbb{E}_{q(z|x)} \left[ \log \frac{p(z|x)p(x)}{q(z|x)} \right] \\
 &= \int q(z|x) \log \frac{p(z|x)}{q(z|x)} dz + \underbrace{\log p(x)}_{\text{constant}} \\
 &= -\text{KL} (q(z|x) \parallel p(z|x)) + \log p(x)
 \end{aligned}$$

# An approximate posterior

$$\begin{aligned}
 \log p(x) &\geq \underbrace{\mathbb{E}_{q(z|x)} \left[ \log \frac{p(x, z)}{q(z|x)} \right]}_{\text{ELBO}} \\
 &= \mathbb{E}_{q(z|x)} \left[ \log \frac{p(z|x)p(x)}{q(z|x)} \right] \\
 &= \int q(z|x) \log \frac{p(z|x)}{q(z|x)} dz + \underbrace{\log p(x)}_{\text{constant}} \\
 &= -\text{KL} (q(z|x) \parallel p(z|x)) + \log p(x)
 \end{aligned}$$

We have derived a lower bound on the log-evidence whose gap is exactly  $\text{KL} (q(z|x) \parallel p(z|x))$ .

# Variational Inference

## Objective

$$\max_{q(z|x)} \mathbb{E} [\log p(x, z)] + \mathbb{H}(q(z|x))$$

- The ELBO is a lower bound on  $\log p(x)$

# Mean field assumption

Suppose we have  $N$  latent variables

- assume the posterior factorises as  $N$  independent terms
- each with an independent set of parameters

$$q(z_1, \dots, z_N) = \underbrace{\prod_{i=1}^N q_{\lambda_i}(z_i)}_{\text{mean field}}$$

# Amortised variational inference

Amortise the cost of inference using NNs

$$q(z_1, \dots, z_N | x_1, \dots, x_N) = \prod_{i=1}^N q_\lambda(z_i | x_i)$$

with a shared set of parameters

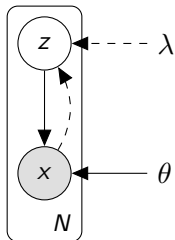
- e.g.  $Z|x \sim \mathcal{N}(\underbrace{\mu_\lambda(x), \sigma_\lambda(x)^2}_{\text{inference network}})$

# Outline

- 1 Deep generative models
- 2 Variational inference
- 3 Variational auto-encoder

# Variational auto-encoder

Generative model with NN likelihood



- complex (non-linear) observation model  $p_{\theta}(x|z)$
- complex (non-linear) mapping from data to latent variables  $q_{\lambda}(z|x)$

Jointly optimise generative model  $p_{\theta}(x|z)$  and inference model  $q_{\lambda}(z|x)$  under the same objective (ELBO)

# Objective

$$\log p(x) \geq \overbrace{\mathbb{E}_{q_\lambda(z|x)} [\log p_\theta(x, z)]}^{\text{ELBO}} + \mathbb{H}(q_\lambda(z|x))$$



# Objective

$$\begin{aligned}\log p(x) &\geq \overbrace{\mathbb{E}_{q_\lambda(z|x)} [\log p_\theta(x, z)] + \mathbb{H}(q_\lambda(z|x))}^{\text{ELBO}} \\ &= \mathbb{E}_{q_\lambda(z|x)} [\log p_\theta(x|z) + \log p(z)] + \mathbb{H}(q_\lambda(z|x))\end{aligned}$$

# Objective

$$\begin{aligned}
 \log p(x) &\geq \overbrace{\mathbb{E}_{q_\lambda(z|x)} [\log p_\theta(x, z)] + \mathbb{H}(q_\lambda(z|x))}^{\text{ELBO}} \\
 &= \mathbb{E}_{q_\lambda(z|x)} [\log p_\theta(x|z) + \log p(z)] + \mathbb{H}(q_\lambda(z|x)) \\
 &= \mathbb{E}_{q_\lambda(z|x)} [\log p_\theta(x|z)] - \text{KL}(q_\lambda(z|x) || p(z))
 \end{aligned}$$

# Objective

$$\begin{aligned}
 \log p(x) &\geq \overbrace{\mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x, z)] + \mathbb{H}(q_{\lambda}(z|x))}^{\text{ELBO}} \\
 &= \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z) + \log p(z)] + \mathbb{H}(q_{\lambda}(z|x)) \\
 &= \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \text{KL}(q_{\lambda}(z|x) || p(z))
 \end{aligned}$$

Parameter estimation

$$\arg \max_{\theta, \lambda} \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \text{KL}(q_{\lambda}(z|x) || p(z))$$

# Objective

$$\begin{aligned}
 \log p(x) &\geq \overbrace{\mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x, z)] + \mathbb{H}(q_{\lambda}(z|x))}^{\text{ELBO}} \\
 &= \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z) + \log p(z)] + \mathbb{H}(q_{\lambda}(z|x)) \\
 &= \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \text{KL}(q_{\lambda}(z|x) \parallel p(z))
 \end{aligned}$$

Parameter estimation

$$\arg \max_{\theta, \lambda} \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \text{KL}(q_{\lambda}(z|x) \parallel p(z))$$

- assume  $\text{KL}(q_{\lambda}(z|x) \parallel p(z))$  analytical  
true for exponential families

# Objective

$$\begin{aligned}
 \log p(x) &\geq \overbrace{\mathbb{E}_{q_\lambda(z|x)} [\log p_\theta(x, z)] + \mathbb{H}(q_\lambda(z|x))}^{\text{ELBO}} \\
 &= \mathbb{E}_{q_\lambda(z|x)} [\log p_\theta(x|z) + \log p(z)] + \mathbb{H}(q_\lambda(z|x)) \\
 &= \mathbb{E}_{q_\lambda(z|x)} [\log p_\theta(x|z)] - \text{KL}(q_\lambda(z|x) \parallel p(z))
 \end{aligned}$$

## Parameter estimation

$$\arg \max_{\theta, \lambda} \mathbb{E}_{q_\lambda(z|x)} [\log p_\theta(x|z)] - \text{KL}(q_\lambda(z|x) \parallel p(z))$$

- assume  $\text{KL}(q_\lambda(z|x) \parallel p(z))$  analytical  
true for exponential families
- approximate  $\mathbb{E}_{q_\lambda(z|x)} [\log p_\theta(x|z)]$  by sampling  
true because we design  $q_\lambda(z|x)$  to be simple

# Generative Network Gradient

$$\frac{\partial}{\partial \theta} \left( \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \overbrace{\text{KL}(q_{\lambda}(z|x) || p(z))}^{\text{constant wrt } \theta} \right)$$

# Generative Network Gradient

$$\begin{aligned} & \frac{\partial}{\partial \theta} \left( \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \overbrace{\text{KL}(q_{\lambda}(z|x) || p(z))}^{\text{constant wrt } \theta} \right) \\ &= \mathbb{E}_{q_{\lambda}(z|x)} \left[ \frac{\partial}{\partial \theta} \log p_{\theta}(x|z) \right] \end{aligned}$$

# Generative Network Gradient

$$\begin{aligned}
 & \frac{\partial}{\partial \theta} \left( \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \overbrace{\text{KL}(q_{\lambda}(z|x) || p(z))}^{\text{constant wrt } \theta} \right) \\
 &= \mathbb{E}_{q_{\lambda}(z|x)} \left[ \frac{\partial}{\partial \theta} \log p_{\theta}(x|z) \right] \\
 &\stackrel{\text{MC}}{\approx} \frac{1}{S} \sum_{i=1}^S \frac{\partial}{\partial \theta} \log p_{\theta}(x|z_i)
 \end{aligned}$$



# Generative Network Gradient

$$\begin{aligned}
 & \frac{\partial}{\partial \theta} \left( \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \overbrace{\text{KL}(q_{\lambda}(z|x) || p(z))}^{\text{constant wrt } \theta} \right) \\
 &= \mathbb{E}_{q_{\lambda}(z|x)} \left[ \frac{\partial}{\partial \theta} \log p_{\theta}(x|z) \right] \\
 &\approx^{\text{MC}} \frac{1}{S} \sum_{i=1}^S \frac{\partial}{\partial \theta} \log p_{\theta}(x|z_i)
 \end{aligned}$$

Note:  $q_{\lambda}(z|x)$  does not depend on  $\theta$ .

# Inference Network Gradient

$$\frac{\partial}{\partial \lambda} \left( \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \overbrace{\text{KL}(q_{\lambda}(z|x) || p(z))}^{\text{analytical}} \right)$$

# Inference Network Gradient

$$\begin{aligned} & \frac{\partial}{\partial \lambda} \left( \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \overbrace{\text{KL}(q_{\lambda}(z|x) || p(z))}^{\text{analytical}} \right) \\ &= \frac{\partial}{\partial \lambda} \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \underbrace{\frac{\partial}{\partial \lambda} \text{KL}(q_{\lambda}(z|x) || p(z))}_{\text{analytical computation}} \end{aligned}$$

# Inference Network Gradient

$$\begin{aligned} & \frac{\partial}{\partial \lambda} \left( \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \overbrace{\text{KL}(q_{\lambda}(z|x) || p(z))}^{\text{analytical}} \right) \\ &= \frac{\partial}{\partial \lambda} \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \underbrace{\frac{\partial}{\partial \lambda} \text{KL}(q_{\lambda}(z|x) || p(z))}_{\text{analytical computation}} \end{aligned}$$

The first term again requires approximation by sampling,  
 but there is a problem

# Inference Network Gradient

$$\begin{aligned}
 & \frac{\partial}{\partial \lambda} \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] \\
 &= \frac{\partial}{\partial \lambda} \int q_{\lambda}(z|x) \log p_{\theta}(x|z) dz \\
 &= \underbrace{\int \frac{\partial}{\partial \lambda} (q_{\lambda}(z|x)) \log p_{\theta}(x|z) dz}_{\text{not an expectation}}
 \end{aligned}$$

# Inference Network Gradient

$$\begin{aligned}
 & \frac{\partial}{\partial \lambda} \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] \\
 &= \frac{\partial}{\partial \lambda} \int q_{\lambda}(z|x) \log p_{\theta}(x|z) dz \\
 &= \underbrace{\int \frac{\partial}{\partial \lambda} (q_{\lambda}(z|x)) \log p_{\theta}(x|z) dz}_{\text{not an expectation}}
 \end{aligned}$$

- MC estimator is non-differentiable: cannot sample first

# Inference Network Gradient

$$\begin{aligned}
 & \frac{\partial}{\partial \lambda} \mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] \\
 &= \frac{\partial}{\partial \lambda} \int q_{\lambda}(z|x) \log p_{\theta}(x|z) dz \\
 &= \underbrace{\int \frac{\partial}{\partial \lambda} (q_{\lambda}(z|x)) \log p_{\theta}(x|z) dz}_{\text{not an expectation}}
 \end{aligned}$$

- MC estimator is non-differentiable: cannot sample first
- Differentiating the expression does not yield an expectation: cannot approximate via MC

# Reparametrisation

Find a transformation  $h : z \mapsto \epsilon$  such that  
 $\epsilon$  does not depend on  $\lambda$

- $h(z, \lambda)$  needs to be invertible
- $h(z, \lambda)$  needs to be differentiable

---

(Kingma and Welling, 2013; Rezende et al., 2014; Titsias and  
Lázaro-Gredilla, 2014)



# Reparametrisation

Find a transformation  $h : z \mapsto \epsilon$  such that  
 $\epsilon$  does not depend on  $\lambda$

- $h(z, \lambda)$  needs to be invertible
- $h(z, \lambda)$  needs to be differentiable

Invertibility implies

- $h(z, \lambda) = \epsilon$
- $h^{-1}(\epsilon, \lambda) = z$

---

(Kingma and Welling, 2013; Rezende et al., 2014; Titsias and  
Lázaro-Gredilla, 2014)

# Gaussian Transformation

If  $Z \sim \mathcal{N}(\mu_\lambda(x), \sigma_\lambda(x)^2)$  then

$$h(z, \lambda) = \frac{z - \mu_\lambda(x)}{\sigma_\lambda(x)} = \epsilon \sim \mathcal{N}(0, 1)$$

$$h^{-1}(\epsilon, \lambda) = \mu_\lambda(x) + \sigma_\lambda(x) \odot \epsilon \quad \epsilon \sim \mathcal{N}(0, 1)$$

$$= \frac{\partial}{\partial \lambda} \int q_{\lambda}(z|x) \log p_{\theta}(x|z) \mathrm{d}z$$

# Inference Network – Reparametrised Gradient

$$= \frac{\partial}{\partial \lambda} \int q_{\lambda}(z|x) \log p_{\theta}(x|z) dz$$

$$= \frac{\partial}{\partial \lambda} \int q(\epsilon) \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) d\epsilon$$

$$\begin{aligned}
 &= \frac{\partial}{\partial \lambda} \int q_{\lambda}(z|x) \log p_{\theta}(x|z) \, dz \\
 &= \frac{\partial}{\partial \lambda} \int \mathbf{q}(\epsilon) \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) \, d\epsilon \\
 &= \int \mathbf{q}(\epsilon) \frac{\partial}{\partial \lambda} \left[ \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) \right] \, d\epsilon
 \end{aligned}$$

# Inference Network – Reparametrised Gradient

$$= \frac{\partial}{\partial \lambda} \int q_{\lambda}(z|x) \log p_{\theta}(x|z) dz$$

$$= \frac{\partial}{\partial \lambda} \int \mathbf{q}(\epsilon) \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) d\epsilon$$

$$= \int q(\epsilon) \frac{\partial}{\partial \lambda} \left[ \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) \right] d\epsilon$$

$$= \int q(\epsilon) \underbrace{\frac{\partial}{\partial z} \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) \times \frac{\partial}{\partial \lambda} h^{-1}(\epsilon, \lambda)}_{\text{chain rule}} d\epsilon$$

# Inference Network – Reparametrised Gradient

$$\begin{aligned}
 &= \frac{\partial}{\partial \lambda} \int q_{\lambda}(z|x) \log p_{\theta}(x|z) \, dz \\
 &= \frac{\partial}{\partial \lambda} \int \mathbf{q}(\epsilon) \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) \, d\epsilon \\
 &= \int q(\epsilon) \frac{\partial}{\partial \lambda} \left[ \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) \right] \, d\epsilon \\
 &= \int q(\epsilon) \underbrace{\frac{\partial}{\partial z} \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) \times \frac{\partial}{\partial \lambda} h^{-1}(\epsilon, \lambda)}_{\text{chain rule}} \, d\epsilon \\
 &= \mathbb{E}_{q(\epsilon)} \left[ \frac{\partial}{\partial z} \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) \times \frac{\partial}{\partial \lambda} h^{-1}(\epsilon, \lambda) \right]
 \end{aligned}$$

## Inference Network – Reparametrised Gradient

$$= \mathbb{E}_{q(\epsilon)} \left[ \frac{\partial}{\partial z} \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) \times \frac{\partial}{\partial \lambda} h^{-1}(\epsilon, \lambda) \right]$$



# Inference Network – Reparametrised Gradient

$$\begin{aligned}
 &= \mathbb{E}_{q(\epsilon)} \left[ \frac{\partial}{\partial z} \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) \times \frac{\partial}{\partial \lambda} h^{-1}(\epsilon, \lambda) \right] \\
 &\stackrel{\text{MC}}{\approx} \frac{1}{S} \sum_{i=1}^S \underbrace{\frac{\partial}{\partial z} \log p_{\theta}(x | \overbrace{h^{-1}(\epsilon, \lambda)}^{=z}) \times \frac{\partial}{\partial \lambda} h^{-1}(\epsilon, \lambda)}_{\text{backprop's job}}
 \end{aligned}$$

# Gaussian KL

## ELBO

$$\mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \text{KL} (q_{\lambda}(z|x) || p(z))$$

# Gaussian KL

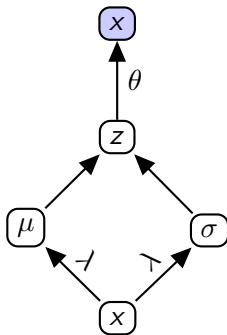
## ELBO

$$\mathbb{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - \text{KL} (q_{\lambda}(z|x) || p(z))$$

Analytical computation of  $-\text{KL} (q_{\lambda}(z|x) || p(z))$ :

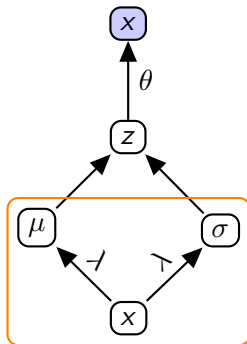
$$\frac{1}{2} \sum_{i=1}^d (1 + \log (\sigma_i^2) - \mu_i^2 - \sigma_i^2)$$

# Computation Graph



# Computation Graph

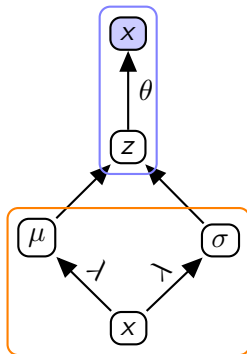
inference model



# Computation Graph

generative model

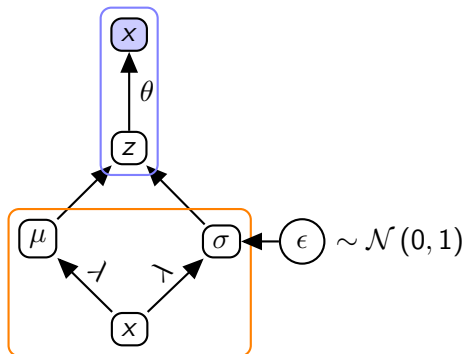
inference model



# Computation Graph

generative model

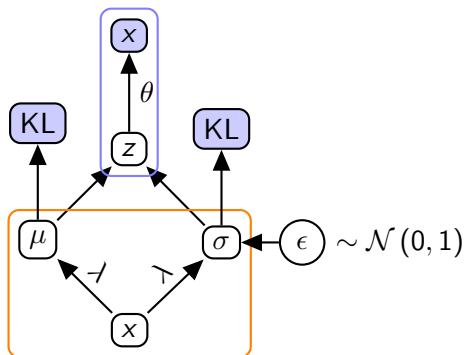
inference model



# Computation Graph

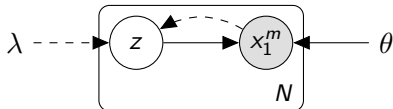
generative model

inference model





## Example



Generative model

- $Z \sim \mathcal{N}(0, I)$
- $X_i | z, x_{<i} \sim \text{Cat}(f_\theta(z, x_{<i}))$

Inference model

- $Z \sim \mathcal{N}(\mu_\lambda(x_1^m), \sigma_\lambda(x_1^m)^2)$

# VAEs – Summary

## Advantages

- Backprop training
- Easy to implement
- Posterior inference possible
- One objective for both NNs

# VAEs – Summary

## Advantages

- Backprop training
- Easy to implement
- Posterior inference possible
- One objective for both NNs

## Drawbacks

- Discrete latent variables are difficult
- Optimisation may be difficult with several latent variables
- Location-scale families only  
but see Ruiz et al. (2016) and Kucukelbir et al. (2017)

# Summary

## Deep learning in NLP

- task-driven feature extraction
- models with more realistic assumptions

## Probabilistic modelling

- better (or at least more explicit) statistical assumptions
- compact models
- semi-supervised learning

# Literature I

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. 2014. URL <http://arxiv.org/abs/1409.0473>.
- David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational inference: A review for statisticians. 01 2016. URL <https://arxiv.org/abs/1601.00670>.
- Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 10–21, Berlin, Germany, August 2016. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/K16-1002>.
- Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. 2013. URL <http://arxiv.org/abs/1312.6114>.

## Literature II

- Alp Kucukelbir, Dustin Tran, Rajesh Ranganath, Andrew Gelman, and David M. Blei. Automatic differentiation variational inference. *Journal of Machine Learning Research*, 18(14):1–45, 2017. URL <http://jmlr.org/papers/v18/16-107.html>.
- Danilo J. Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, pages 1278–1286, 2014. URL <http://jmlr.org/proceedings/papers/v32/rezende14.pdf>.
- Francisco R Ruiz, Michalis Titsias RC AUEB, and David Blei. The generalized reparameterization gradient. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *NIPS*, pages 460–468. 2016. URL <http://papers.nips.cc/paper/6328-the-generalized-reparameterization-gradient.pdf>.

## Literature III

Rico Sennrich, Barry Haddow, and Alexandra Birch. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany, August 2016. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/P16-1009>.

Michalis Titsias and Miguel Lázaro-Gredilla. Doubly stochastic variational bayes for non-conjugate inference. In Tony Jebara and Eric P. Xing, editors, *ICML*, pages 1971–1979, 2014. URL <http://jmlr.org/proceedings/papers/v32/titsias14.pdf>.