# Variational Inference
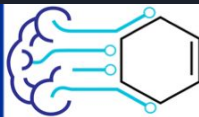## From basics to modern applications

Adam Arany

17 October 2022
3rd AIDD School, Leuven, Belgium

# Overview

- Bayesian models and inference recap (from Lugano)
  - Bayesian probability
  - Graphical model notation
  - Predictive inference
- Variational Inference Basics
  - Calculus of Variation
  - Derivation of the Evidence Lower Bound (ELBO)
  - Mean Field Approximation
  - Gradient based VI
- Variational Autoencoders
  - Amortized inference
  - Reparametrization trick
- Bayesian NNs

# Bayesian Probability

*Probability of the hypothesis*
*after observing the evidence*
*"a posteriori"*

*Probability of the hypothesis*
*before observing the evidence*
*"a priori"*

posterior     likelihood     prior

$$P(H|D) = \frac{P(D|H)P(H)}{P(D)}$$

marginal
likelihood

$$P(D) = \int P(D|H)P(H)\,dH$$

*This integral is most often intractable!*

- have an epistemic / subjectivist interpretation
- true independent of interpretations

Pierre-Simon Laplace

# Elements of probabilistic models

(x) Random variables
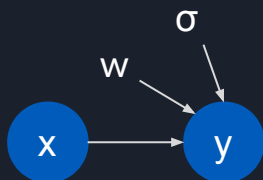$$X \sim \mathcal{N}(0,1)$$
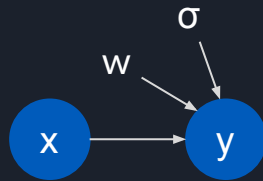
(x) → (y) Dependencies
$$Y \sim \mathcal{N}(2X, 0.5)$$

σ
w
(x) → (y) Parameters
$$Y \sim \mathcal{N}(wX, \sigma)$$

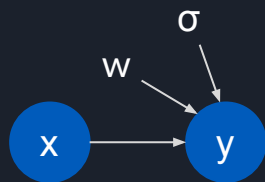Probabilistic graphical models (PGMs), Bayesian networks WARNING!

# Bayesian models

Bayesian treatment:
- Parameter is just a random variable
- We do not expect to find 'the real' parameter value exactly
- We search for the distribution of the parameters supported by the data.



point
parametrization

Bayesian

# Bayesian models

point
parametrization

hyperparameters

$\sigma_w$   $\mu_w$   $\alpha$   $\beta$   $\sigma$   $w$   $x$   $y$

Bayesian

hyperparameters

$\sigma_w$   $\mu_w$   $\alpha$   $\beta$   $\sigma$   $w$   $x$   $y$

hierarchical
Bayesian

# Frequently used shorthand notations

1) Vector, matrix, tensor valued random variables

$$y \sim \mathcal{N}(W^{\top}x, \Sigma)$$

2) Plate models

# Predictive inference

x   Unobserved variable

x   Observed variable

Predicted outcome? Just another random variable

$$P(y'|x',x,y) = P(y'|x',w,\sigma)P(w,\sigma|x,y)$$

$\underbrace{\phantom{(x,y)}}_{\mathcal{D}}$    $\underbrace{\phantom{(x,y)}}_{\mathcal{D}}$

Posterior predictive distribution

(model) posterior

What is the mean?

$$\mathbb{E}[y'|x',\mathcal{D}] = \mathbb{E}_{P(w,\sigma|\mathcal{D})}\big[f_{w,\sigma}(x')\big]$$

"Bayesian model averaging"

Note the similarity with ensembles!

# Calculus of variation



**Calculus:** concerned with functions $f(x)$ mapping "numbers" to "numbers"

**Variational calculus:** concerned with $F[f]$ functionals mapping functions to "numbers"

Searching for extremal points of F functional

Several physics problem:
- Brachistochrone problem
- Shape of hanging chain
- Soap bubbles
- Fermat principle
- Principle of least action / principle of stationary action

# Calculus of variation

$$H[P] = \int_{-\infty}^{\infty} P(x)\log P(x)dx$$

Entropy

$$\mathcal{S}[L] = \int_{t_1}^{t_2} L(q(t),\dot{q}(t))dt$$

Action

$$D[Q] = D_{KL}(Q\|P) = \int_{-\infty}^{\infty} Q(x)\log\frac{Q(x)}{P(x)}dx$$

Kullback-Leibler divergence

# Calculus of variation

$$H[P] = \int_{-\infty}^{\infty} P(x) \log P(x) \, dx$$

Entropy

$$S[L] = \int_{t_1}^{t_2} L(q(t), \dot{q}(t)) \, dt$$

Action

$$D[Q] = D_{KL}(Q \| P) = \int_{-\infty}^{\infty} Q(x) \log \frac{Q(x)}{P(x)} \, dx$$

Kullback-Leibler divergence

# Calculus of variation

Searching for a distribution (a function) Q, that minimize a KL-divergence (a functional) is a problem in the calculus of variation.

Therefore the inference method using this approximation is named Variational.

$$D[Q] = D_{KL}(Q\|P) = \int_{-\infty}^{\infty} Q(x) \log \frac{Q(x)}{P(x)} dx$$

Kullback-Leibler divergence

# General inference problem

Given a model (here as a graphical model).
Define a set of observed variables E

$$E = \{x, y\}$$

And the variables of interest H

$$H = \{w, \sigma, \sigma_w\}$$

We want to estimate the posterior of H conditioned on E, P(H|E)

$$p(w, \sigma, \sigma_w | x, y)$$

# Variational inference

If we would have $p(w, \sigma, \sigma_w | x, y)$
in analytical form, our job would be done.
- Often there is no such analytical form
- Search for a function $q(w, \sigma, \sigma_w) \approx p(w, \sigma, \sigma_w | x, y)$

$$\min_{\phi} D_{KL}(q_{\phi}(w, \sigma, \sigma_w) \| p(w, \sigma, \sigma_w | x, y))$$

$\phi$: variational parameter

# Variational inference

If we would have $p(w, \sigma, \sigma_w | x, y)$
in analytical form, our job would be done.

- Often there is no such analytical form
- Search for a function $q(w, \sigma, \sigma_w) \approx p(w, \sigma, \sigma_w | x, y)$

$$\min_{\phi} D_{KL}(q_{\phi}(w, \sigma, \sigma_w) \| p(w, \sigma, \sigma_w | x, y))$$

$\phi$ : variational parameter

Every evidence result in a different q!

Every inference query require optimization. This is a price we pay.

# Variational inference

If we would have $p(w, \sigma, \sigma_w | x, y)$
in analytical form, our job would be done.
- Often there is no such analytical form
- Search for a function $q(w, \sigma, \sigma_w) \approx p(w, \sigma, \sigma_w | x, y)$

$$\min_{\phi} D_{KL}(q_\phi(w, \sigma, \sigma_w) \| p(w, \sigma, \sigma_w | x, y))$$

$\phi$: variational parameter



It can be shown that equivalently we can take the following objective:

$$\min_{\phi} \underbrace{D_{KL}(q_\phi(w, \sigma, \sigma_w) \| p(w, \sigma, \sigma_w)) - \mathbb{E}_{q_\phi(w, \sigma, \sigma_w)}[\log p(x, y | w, \sigma, \sigma_w)]}$$
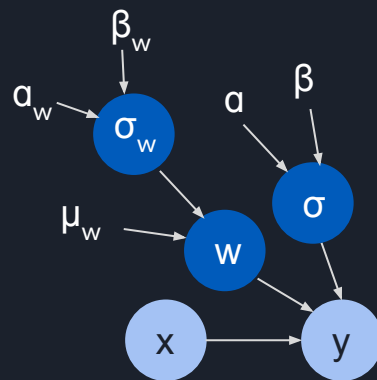
Evidence lower bound (ELBO)

# Variational inference

If we would have $p(w,\sigma,\sigma_w|x,y)$
in analytical form, our job would be done.
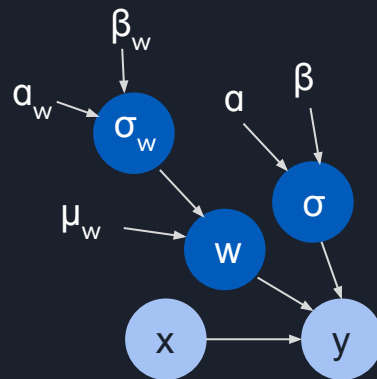
- Often there is no such analytical form
- Search for a function $q(w,\sigma,\sigma_w) \approx p(w,\sigma,\sigma_w|x,y)$

$$\min_\phi D_{KL}(q_\phi(w,\sigma,\sigma_w) \| p(w,\sigma,\sigma_w|x,y))$$

$\phi$ : variational parameter



It can be shown that equivalently we can take the following objective:

$$\min_\phi \underbrace{D_{KL}(q_\phi(w,\sigma,\sigma_w) \| p(w,\sigma,\sigma_w))}_{\text{Regularization}} - \underbrace{\mathbb{E}_{q_\phi(w,\sigma,\sigma_w)}[\log p(x,y|w,\sigma,\sigma_w)]}_{\text{Expected likelihood}}$$

# Derivation of the ELBO

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H|E))$$

# Derivation of the ELBO

$$\min_{\phi} D_{KL}(q_\phi(H)\|p(H|E))$$

$$\min_{\phi}\int q_\phi(H)\log\frac{q_\phi(H)}{p(H|E)}dH=\mathbb{E}_{q_\phi(H)}\big[\log q_\phi(H)-\log p(H|E)\big]$$

# Derivation of the ELBO

$$\min_{\phi} D_{KL}(q_\phi(H) \| p(H|E))$$

$$\min_{\phi} \int q_\phi(H) \log \frac{q_\phi(H)}{p(H|E)} dH = \mathbb{E}_{q_\phi(H)}\left[\log q_\phi(H) - \log p(H|E)\right]$$

$$\frac{p(E|H)p(H)}{P(E)}$$

# Derivation of the ELBO

$$\frac{p(E|H)p(H)}{P(E)}$$

$$\min_{\phi} D_{KL}(q_\phi(H)\|p(H|E))$$

$$\min_{\phi} \int q_\phi(H)\log\frac{q_\phi(H)}{p(H|E)}dH = \mathbb{E}_{q_\phi(H)}\big[\log q_\phi(H) - \log p(H|E)\big]$$

$$\min_{\phi} \mathbb{E}_{q_\phi(H)}\big[\log q_\phi(H) - \log p(H) - \log p(E|H) + \log p(E)\big]$$

# Derivation of the ELBO

$$\frac{p(E|H)p(H)}{P(E)}$$

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H|E))$$

$$\min_{\phi} \int q_{\phi}(H)\log\frac{q_{\phi}(H)}{p(H|E)}dH = \mathbb{E}_{q_{\phi}(H)}\left[\log q_{\phi}(H) - \log p(H|E)\right]$$

$$\min_{\phi} \mathbb{E}_{q_{\phi}(H)}\left[\log q_{\phi}(H) - \log p(H) - \log p(E|H) + \log p(E)\right]$$

# Derivation of the ELBO

$$\frac{p(E|H)p(H)}{P(E)}$$

$$\min_{\phi} D_{KL}(q_\phi(H) \| p(H|E))$$

$$\min_{\phi} \int q_\phi(H) \log \frac{q_\phi(H)}{p(H|E)} dH = \mathbb{E}_{q_\phi(H)} \left[ \log q_\phi(H) - \log \boxed{p(H|E)} \right]$$

$$\min_{\phi} \mathbb{E}_{q_\phi(H)} \left[ \log q_\phi(H) - \log p(H) - \log p(E|H) + \log p(E) \right]$$

# Derivation of the ELBO

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H|E))$$

$$\min_{\phi} \int q_{\phi}(H)\log\frac{q_{\phi}(H)}{p(H|E)}dH = \mathbb{E}_{q_{\phi}(H)}\Big[\log q_{\phi}(H) - \log p(H|E)\Big]$$

$$\min_{\phi} \mathbb{E}_{q_{\phi}(H)}\Big[\log q_{\phi}(H) - \log p(H) - \log p(E|H) + \log p(E)\Big]$$

$$\frac{p(E|H)p(H)}{P(E)}$$

# Derivation of the ELBO

$$\frac{p(E|H)p(H)}{P(E)}$$

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H|E))$$

$$\min_{\phi} \int q_{\phi}(H)\log\frac{q_{\phi}(H)}{p(H|E)}dH = \mathbb{E}_{q_{\phi}(H)}\Big[\log q_{\phi}(H) - \log p(H|E)\Big]$$

$$\min_{\phi} \mathbb{E}_{q_{\phi}(H)}\Big[\log q_{\phi}(H) - \log p(H) - \log p(E|H) + \log p(E)\Big]$$

# Derivation of the ELBO

$$\frac{p(E|H)p(H)}{P(E)}$$

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H|E))$$

$$\min_{\phi} \int q_{\phi}(H)\log\frac{q_{\phi}(H)}{p(H|E)}dH = \mathbb{E}_{q_{\phi}(H)}\Big[\log q_{\phi}(H) - \log \boxed{p(H|E)}\Big]$$

$$\min_{\phi} \mathbb{E}_{q_{\phi}(H)}\Big[\underbrace{\log q_{\phi}(H) - \log p(H)}_{D_{KL}(q_{\phi}(H)\|p(H))} - \log p(E|H) + \log p(E)\Big]$$

# Derivation of the ELBO

$$\frac{p(E|H)p(H)}{P(E)}$$

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H|E))$$

$$\min_{\phi} \int q_{\phi}(H)\log\frac{q_{\phi}(H)}{p(H|E)}dH = \mathbb{E}_{q_{\phi}(H)}\Big[\log q_{\phi}(H) - \log \boxed{p(H|E)}\Big]$$

$$\min_{\phi} \mathbb{E}_{q_{\phi}(H)}\Big[\underbrace{\log q_{\phi}(H) - \log p(H)}_{D_{KL}(q_{\phi}(H)\|p(H))} - \log p(E|H) + \log p(E)\Big]$$

# Derivation of the ELBO

$$\frac{p(E|H)p(H)}{P(E)}$$

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H|E))$$

$$\min_{\phi} \int q_{\phi}(H)\log\frac{q_{\phi}(H)}{p(H|E)}dH = \mathbb{E}_{q_{\phi}(H)}\Big[\log q_{\phi}(H) - \log \boxed{p(H|E)}\Big]$$

$$\min_{\phi} \mathbb{E}_{q_{\phi}(H)}\Big[\underbrace{\log q_{\phi}(H) - \log p(H)}_{D_{KL}(q_{\phi}(H)\|p(H))} - \log p(E|H) + \log p(E)\Big]$$

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H)) - \mathbb{E}_{q_{\phi}(H)}\Big[\log p(E|H)\Big] + \text{const.}$$

# Variational Autoencoders

Diederik P. Kingma and Max Welling. "Auto-encoding variational bayes"

Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. "Stochastic backpropagation and variational inference in deep latent gaussian models"

**KU LEUVEN**

# (Classical) Autoencoders



$$f(.) \qquad g(.)$$

$$L(g(f(x)), x)$$

# Variational Autoencoder



$p_\theta(z)$ prior e.g.: $\mathcal{N}(0, \boldsymbol{I})$

$p_\theta(x|z)$ likelihood e.g: $\mathcal{N}(g_\mu(z), g_\sigma(z))$

# Variational Autoencoder



$p_\theta(z)$ prior e.g.: $\mathcal{N}(0, \boldsymbol{I})$

$p_\theta(x|z)$ likelihood e.g: $\mathcal{N}(g_\mu(z), g_\sigma(z))$

In real life applications often the mean is regarded as output

1, .., N

# Variational Autoencoder



$p_\theta(z)$ prior e.g.: $\mathcal{N}(0, \boldsymbol{I})$

$p_\theta(x|z)$ likelihood e.g: $\mathcal{N}(g_\mu(z), g_\sigma(z))$

In real life applications often the mean is regarded as output

$\phi$

$z$

$\theta$

$z \longrightarrow g_\theta \longrightarrow \mu_x$
$\qquad\qquad \longrightarrow \sigma_x$

$x \longrightarrow f_\phi \longrightarrow \mu_z$
$\qquad\qquad \longrightarrow \sigma_z$

1, .., N

Search for a variational distribution in the form $\mathcal{N}(f_\mu(x), f_\sigma(x))$
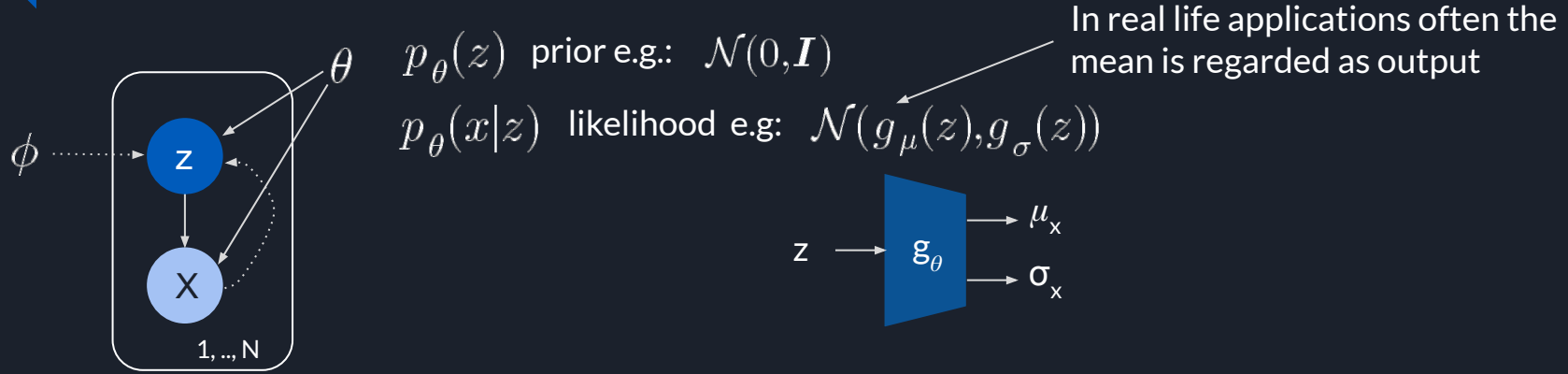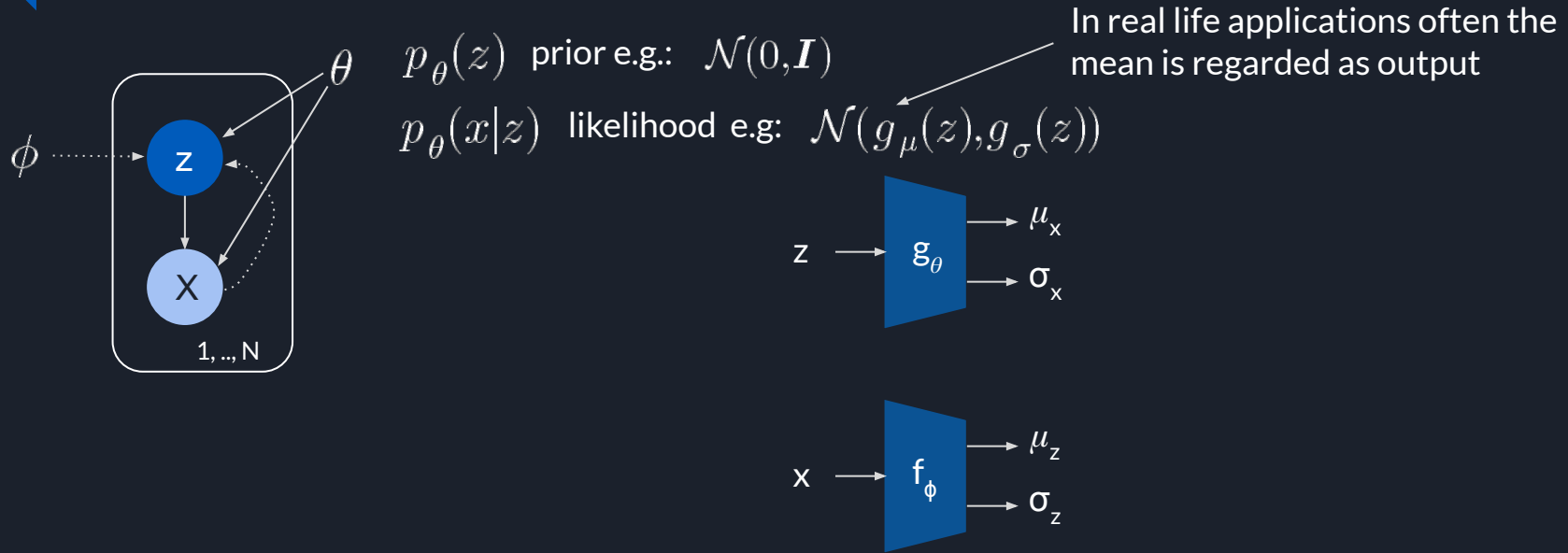Note: an entire family of variational distributions is learned!

# Variational Autoencoder

$p_\theta(z)$ prior e.g.: $\mathcal{N}(0, \boldsymbol{I})$

$p_\theta(x|z)$ likelihood e.g: $\mathcal{N}(g_\mu(z), g_\sigma(z))$

In real life applications often the mean is regarded as output



generative model

recognition model

SAMPLING
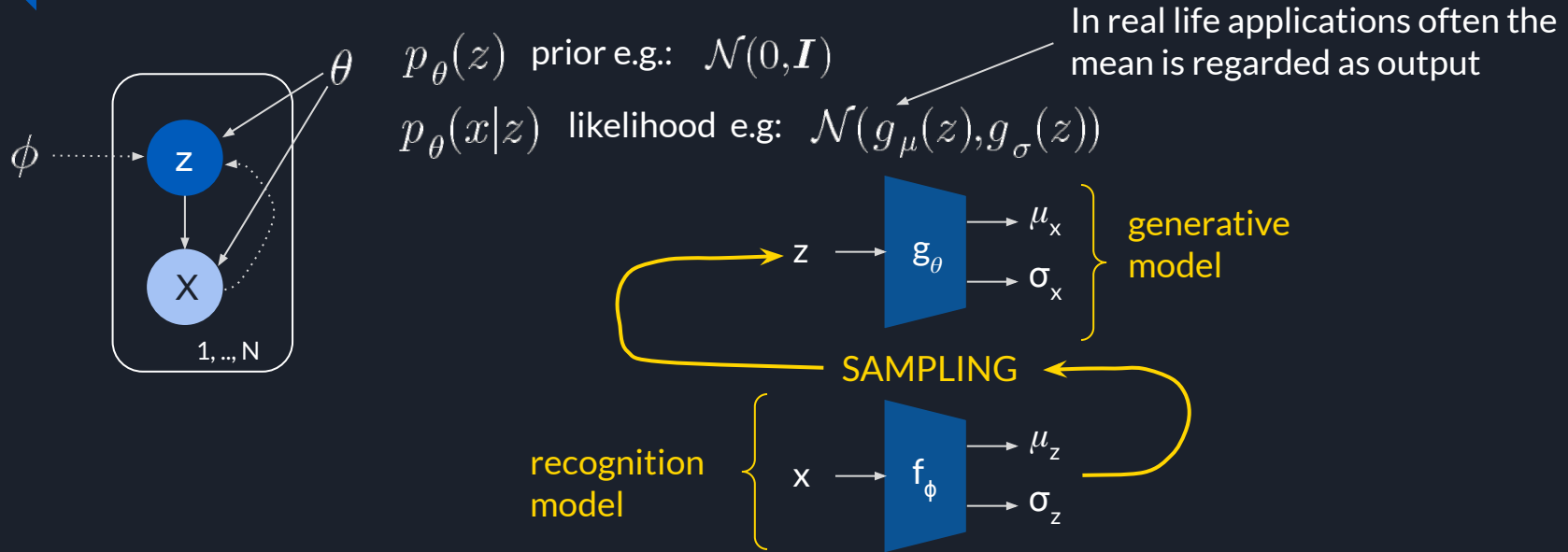
Search for a variational distribution in the form $\mathcal{N}(f_\mu(x), f_\sigma(x))$
Note: an entire family of variational distributions is learned!

# Variational Autoencoder



$p_\theta(z)$ prior e.g.: $\mathcal{N}(0, \boldsymbol{I})$

$p_\theta(x|z)$ likelihood e.g: $\mathcal{N}(g_\mu(z), g_\sigma(z))$

In real life applications often the mean is regarded as output

generative model

SAMPLING

recognition model

Search for a variational distribution in the form $\mathcal{N}(f_\mu(x), f_\sigma(x))$
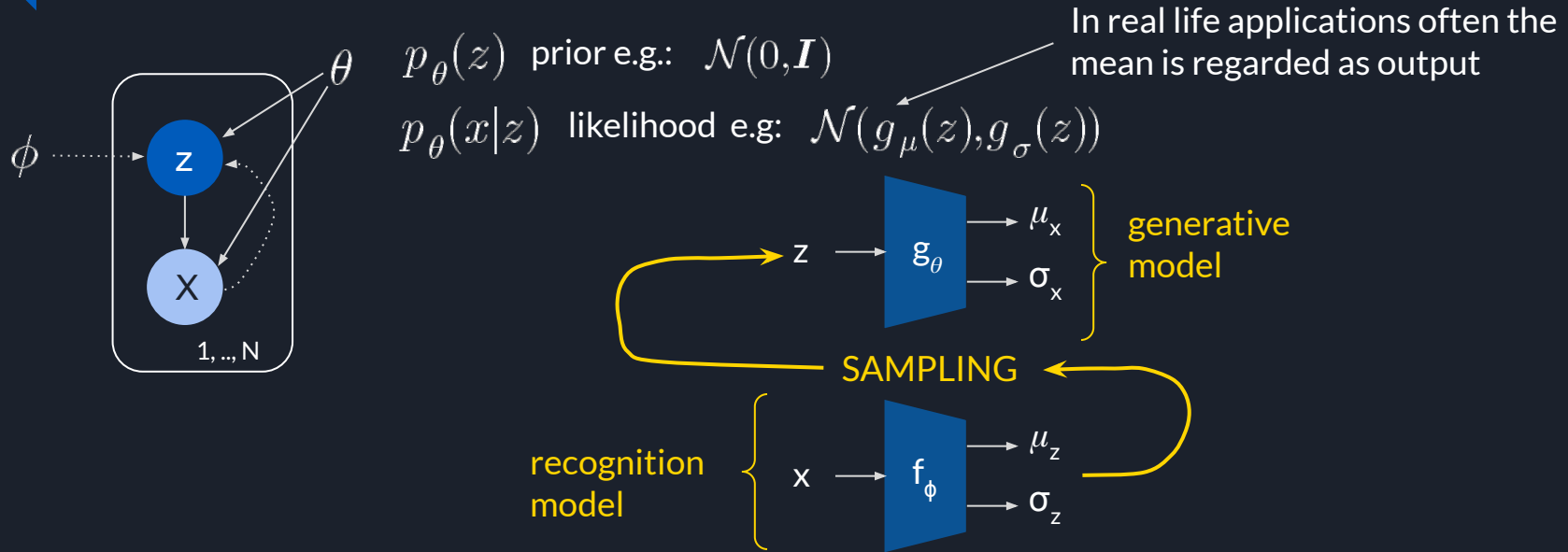Note: an entire family of variational distributions is learned!

Learning $\theta$ and $\phi$ jointly.

# Reparameterization trick

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H)) - \mathbb{E}_{q_{\phi}(H)}\big[\log p(E|H)\big] + \text{const.}$$

We need to evaluate gradients of the following form:     $\nabla_{\phi} \mathbb{E}_{q_{\phi}(H)}\big[f(H)\big]$

# Reparameterization trick

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H)) - \mathbb{E}_{q_{\phi}(H)}\big[\log p(E|H)\big] + \text{const.}$$

We need to evaluate gradients of the following form:　　$\nabla_{\phi}\, \mathbb{E}_{q_{\phi}(H)}\big[f(H)\big]$

Let us assume the following form:　$q_{\phi}(H) = \mathcal{N}(\mu(X), \sigma(X))$

# Reparameterization trick

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H)) - \mathbb{E}_{q_{\phi}(H)}\big[\log p(E|H)\big] + \text{const.}$$

We need to evaluate gradients of the following form: $\quad \nabla_{\phi}\, \mathbb{E}_{q_{\phi}(H)}\big[f(H)\big]$

Let us assume the following form: $\quad q_{\phi}(H) = \mathcal{N}(\mu(X), \sigma(X))$

Samples from the distribution can be generated as:

$$\mu(X) + \epsilon\sigma(X); \ \epsilon \sim \mathcal{N}(0, I)$$

# Reparameterization trick

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H)) - \mathbb{E}_{q_{\phi}(H)}\big[\log p(E|H)\big] + \text{const.}$$

We need to evaluate gradients of the following form: $\quad \nabla_{\phi}\,\mathbb{E}_{q_{\phi}(H)}\big[f(H)\big]$

Let us assume the following form: $q_{\phi}(H) = \mathcal{N}(\mu(X),\sigma(X))$

Samples from the distribution can be generated as:

$$\mu(X) + \epsilon\sigma(X); \; \epsilon \sim \mathcal{N}(0,I)$$

Substituting to the gradient expression, we get the reparameterized form:

$$\nabla_{\phi}\mathbb{E}_{\mathcal{N}(\epsilon|0,I)}\big[f(\mu(X)+\epsilon\sigma(X))\big]$$

# Reparameterization trick

$$\min_{\phi} D_{KL}(q_{\phi}(H)\|p(H)) - \mathbb{E}_{q_{\phi}(H)}\big[\log p(E|H)\big] + \text{const.}$$

We need to evaluate gradients of the following form: $\nabla_{\phi} \mathbb{E}_{q_{\phi}(H)}\big[f(H)\big]$

Let us assume the following form: $q_{\phi}(H) = \mathcal{N}(\mu(X), \sigma(X))$

Samples from the distribution can be generated as:

$$\mu(X) + \epsilon\sigma(X); \ \epsilon \sim \mathcal{N}(0, I)$$

Substituting to the gradient expression, we get the reparameterized form:

$$\nabla_{\phi} \mathbb{E}_{\underbrace{\mathcal{N}(\epsilon|0,I)}}\big[f(\underbrace{\mu(X) + \epsilon\sigma(X)})\big]$$

parameter free
distribution

function
parameterized
by φ

# Reparameterization trick

$$\min_{\phi} D_{KL}(q_{\phi}(H) \| p(H)) - \mathbb{E}_{q_{\phi}(H)}\big[\log p(E|H)\big] + \text{const.}$$

We need to evaluate gradients of the following form: $\nabla_{\phi} \mathbb{E}_{q_{\phi}(H)}\big[f(H)\big]$

Let us assume the following form: $q_{\phi}(H) = \mathcal{N}(\mu(X), \sigma(X))$

Samples from the distribution can be generated as:

$$\mu(X) + \epsilon\sigma(X); \ \epsilon \sim \mathcal{N}(0, I)$$

Substituting to the gradient expression, we get the reparameterized form:

$$\nabla_{\phi} \mathbb{E}_{\underbrace{\mathcal{N}(\epsilon|0,I)}}\big[f(\underbrace{\mu(X) + \epsilon\sigma(X)})\big] \approx \frac{1}{L}\sum_{l=1}^{L}{}_{...}$$

parameter free
distribution

function
parameterized
by φ

Approximated by a finite sum
Kingma and Welling used L=1 ;)

# Reparameterization trick

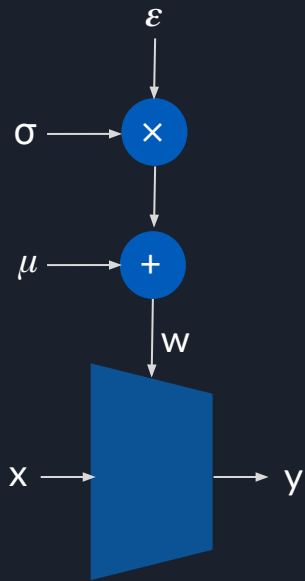$$\ldots = \ldots \left( \ldots (\text{}) \| \ldots (\text{}) \right) = \ldots \left[ \ldots \ldots (\ldots I) \right] + \text{const.}$$

We need to

Let us assu

Samples fr

Substitutin

|  | Target | Basis | Differentiable transformation |
|---|---|---|---|
|  | $\mathcal{N}(\mu, \sigma^2)$ | $\mathcal{N}(0,1)$ | $\mu + \epsilon \sigma$ |
|  | $\mathcal{N}(\mu, RR^\top)$ | $\mathcal{N}(0,I)$ | $\mu + R\epsilon$ |
|  | **any with tractable inverse CDF** | $\mathcal{U}(0,1)$ | **the inverse CDF** |
|  | $\mathcal{C}at(\pi_1, \ldots \pi_k)$ | $\mathcal{G}umbel(0,1)$ | $\underset{i}{\arg\max} \, \epsilon_i + \log \pi_i$ |
|  | … | … | … |

**parameter free distribution**

**function parameterized by φ**

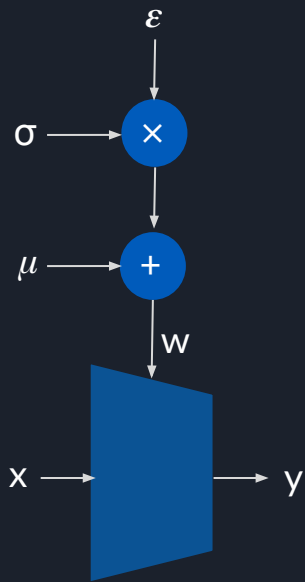**Approximated by a finite sum**
**Kingma and Welling used L=1 ;)**
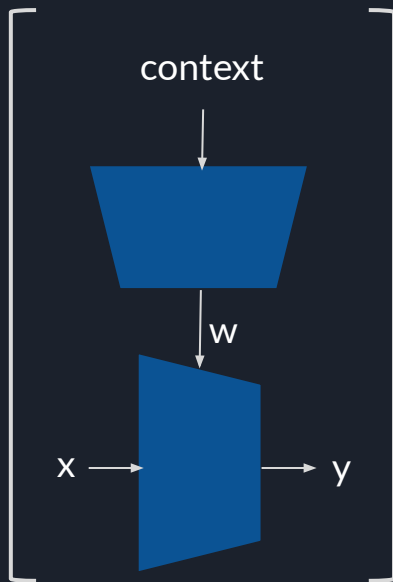
# Bayesian treatment of NN weights



reparameterization trick

# Bayesian treatment of NN weights
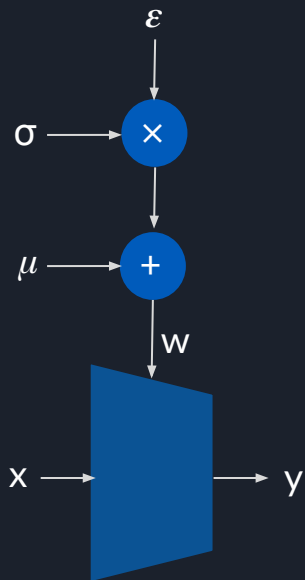


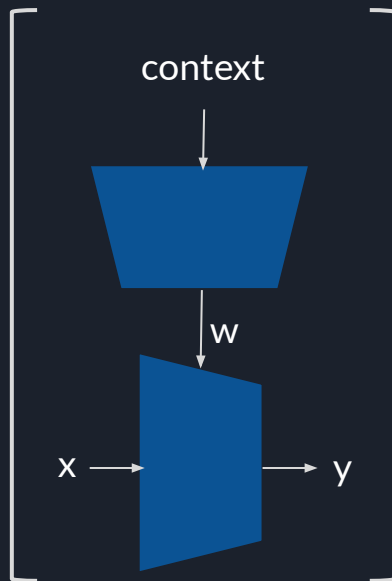reparameterization trick

meta networks

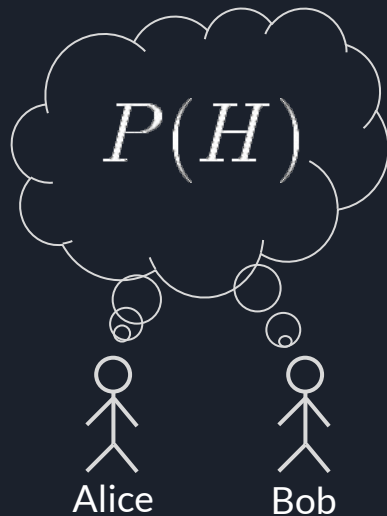# Bayesian treatment of NN weights
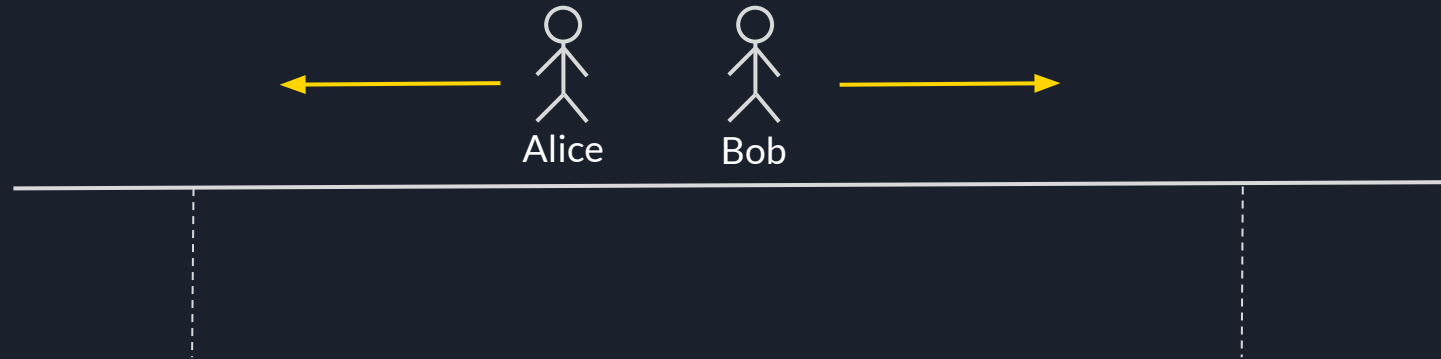

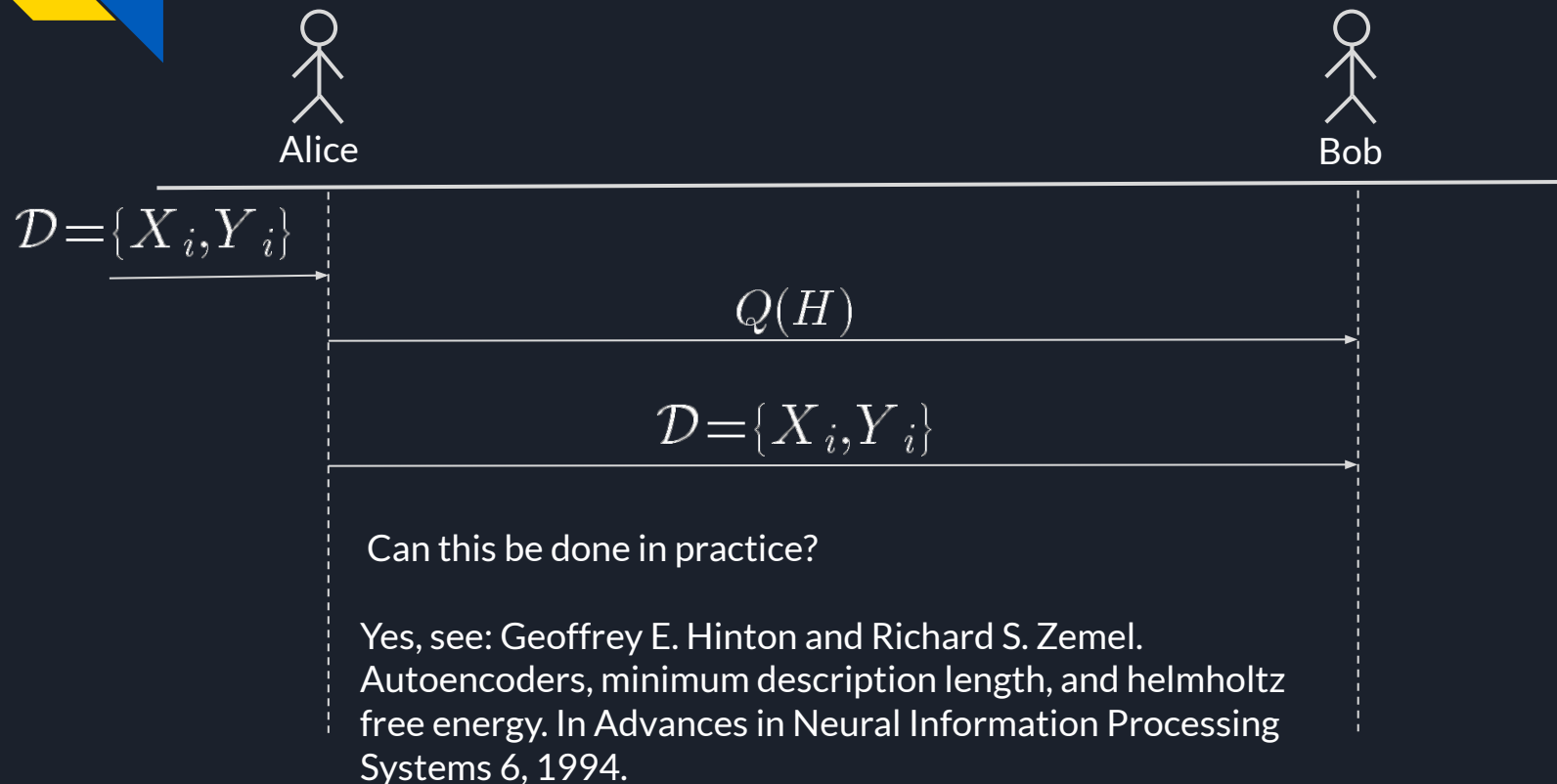
reparameterization trick      meta networks      Bayesian meta networks

# Information theoretical interpretation

# Information theoretical interpretation

# Information theoretical interpretation



Alice

Bob

$\mathcal{D} = \{X_i, Y_i\}$

$Q(H)$

$\mathcal{D} = \{X_i, Y_i\}$

Can this be done in practice?

Yes, see: Geoffrey E. Hinton and Richard S. Zemel. Autoencoders, minimum description length, and helmholtz free energy. In Advances in Neural Information Processing Systems 6, 1994.

Thank you for your attention!