

Auto-Encoding Variational Bayes

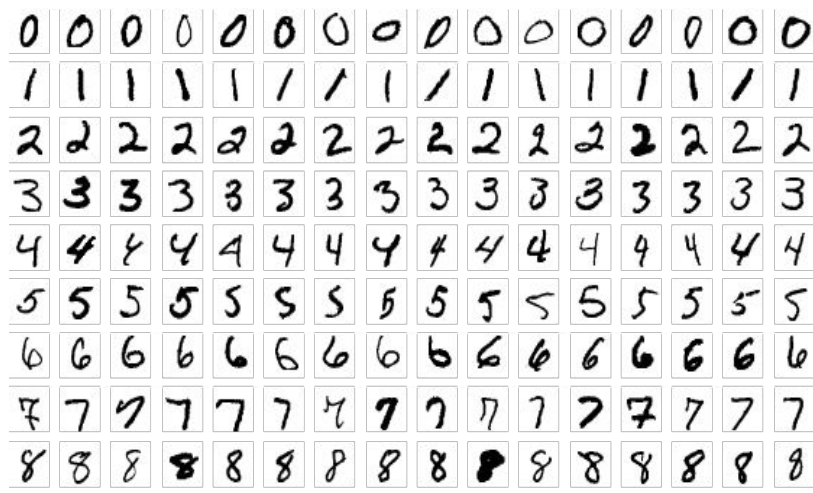
Diederik P. Kingma, Max Welling

Paper report by:

Nikita Ivlev, Lev Leontev, Aleksandr Kariakin

Marginal likelihood

$$\begin{pmatrix} 0.3 \\ 0.5 \\ 0.2 \end{pmatrix} \mathbf{z} \longrightarrow \mathbf{X}$$



Given a set of **independent identically distributed** data points $\mathbf{X} = (x_1, \dots, x_n)$, where $x_i \sim p(x|\theta)$ according to some **probability distribution** parameterized by θ , where θ itself is a **random variable** described by a distribution, i.e. $\theta \sim p(\theta | \alpha)$, the marginal likelihood in general asks what the probability $p(\mathbf{X} | \alpha)$ is, where θ has been **marginalized out** (integrated out):

$$p(\mathbf{X} | \alpha) = \int_{\theta} p(\mathbf{X} | \theta) p(\theta | \alpha) d\theta$$

In our case:

$$\int p_{\theta}(\mathbf{z}) p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z}$$

Problems with existing methods

1. *Intractability*: the case where the integral of the marginal likelihood $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z}$ is intractable (so we cannot evaluate or differentiate the marginal likelihood), where the true posterior density $p_{\theta}(\mathbf{z}|\mathbf{x}) = p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})/p_{\theta}(\mathbf{x})$ is intractable (so the EM algorithm cannot be used), and where the required integrals for any reasonable mean-field VB algorithm are also intractable. These intractabilities are quite common and appear in cases of moderately complicated likelihood functions $p_{\theta}(\mathbf{x}|\mathbf{z})$, e.g. a neural network with a nonlinear hidden layer.
2. *A large dataset*: we have so much data that batch optimization is too costly; we would like to make parameter updates using small minibatches or even single datapoints. Sampling-based solutions, e.g. Monte Carlo EM, would in general be too slow, since it involves a typically expensive sampling loop per datapoint.

Method

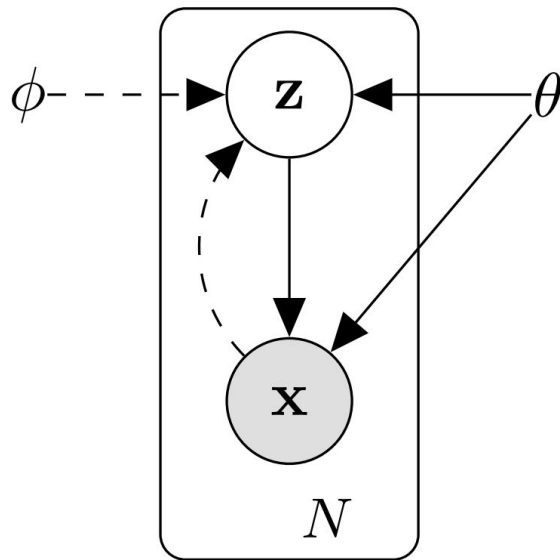


Figure 1: The type of directed graphical model under consideration. Solid lines denote the generative model $p_{\theta}(z)p_{\theta}(x|z)$, dashed lines denote the variational approximation $q_{\phi}(z|x)$ to the intractable posterior $p_{\theta}(z|x)$. The variational parameters ϕ are learned jointly with the generative model parameters θ .

Probabilistic encoder definition

$$p_{\theta}(z|x) \longrightarrow q_{\phi}(z|x)$$

Probabilistic decoder definition

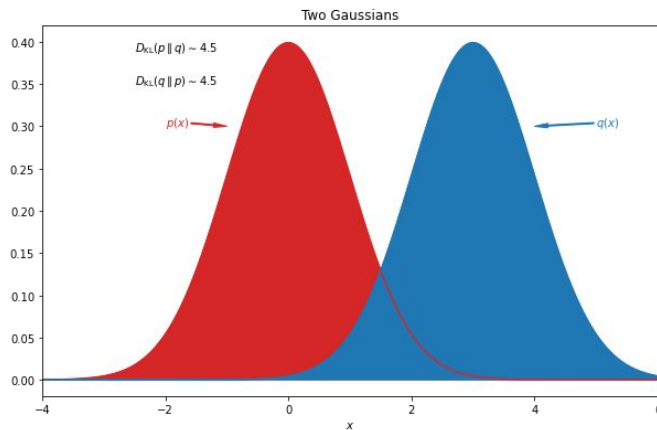
$$p_{\theta}(x|z)$$

This solves intractability!

Kullback–Leibler divergence

The Kullback-Leibler divergence is a measure of the dissimilarity between two probability distributions.

$$D_{\text{KL}}(p \parallel q) + H(p) = H(p, q)$$



Variational lower bound

marginal likelihood

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) = D_{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}^{(i)})) + \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)})$$

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \geq \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)})$$

Now we can optimize $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)})$ instead!

But the gradient of L can't be approximated well with Monte-Carlo methods



Idea of reparametrization

$$\tilde{\mathbf{z}} \sim q_{\phi}(\mathbf{z}|\mathbf{x}) \Rightarrow \tilde{\mathbf{z}} = g_{\phi}(\boldsymbol{\epsilon}, \mathbf{x}) \quad \text{with} \quad \boldsymbol{\epsilon} \sim p(\boldsymbol{\epsilon})$$

1. Tractable inverse CDF. In this case, let $\boldsymbol{\epsilon} \sim \mathcal{U}(\mathbf{0}, \mathbf{I})$, and let $g_{\phi}(\boldsymbol{\epsilon}, \mathbf{x})$ be the inverse CDF of $q_{\phi}(\mathbf{z}|\mathbf{x})$. Examples: Exponential, Cauchy, Logistic, Rayleigh, Pareto, Weibull, Reciprocal, Gompertz, Gumbel and Erlang distributions.
2. Analogous to the Gaussian example, for any "location-scale" family of distributions we can choose the standard distribution (with location = 0, scale = 1) as the auxiliary variable $\boldsymbol{\epsilon}$, and let $g(\cdot) = \text{location} + \text{scale} \cdot \boldsymbol{\epsilon}$. Examples: Laplace, Elliptical, Student's t, Logistic, Uniform, Triangular and Gaussian distributions.
3. Composition: It is often possible to express random variables as different transformations of auxiliary variables. Examples: Log-Normal (exponentiation of normally distributed variable), Gamma (a sum over exponentially distributed variables), Dirichlet (weighted sum of Gamma variates), Beta, Chi-Squared, and F distributions.

Stochastic Gradient Variational Bayes estimator

We can now form Monte Carlo estimates of expectations of some function $f(\mathbf{z})$ w.r.t. $q_\phi(\mathbf{z}|\mathbf{x})$ as follows:

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [f(\mathbf{z})] = \mathbb{E}_{p(\epsilon)} \left[f(g_\phi(\epsilon, \mathbf{x}^{(i)})) \right] \simeq \frac{1}{L} \sum_{l=1}^L f(g_\phi(\epsilon^{(l)}, \mathbf{x}^{(i)})) \quad \text{where} \quad \epsilon^{(l)} \sim p(\epsilon) \quad (5)$$

We apply this technique to the variational lower bound (eq. (2)), yielding our generic Stochastic Gradient Variational Bayes (SGVB) estimator $\tilde{\mathcal{L}}^A(\theta, \phi; \mathbf{x}^{(i)}) \simeq \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$:

$$\tilde{\mathcal{L}}^A(\theta, \phi; \mathbf{x}^{(i)}) = \frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}^{(i)}, \mathbf{z}^{(i,l)}) - \log q_\phi(\mathbf{z}^{(i,l)}|\mathbf{x}^{(i)})$$

where $\mathbf{z}^{(i,l)} = g_\phi(\epsilon^{(i,l)}, \mathbf{x}^{(i)})$ and $\epsilon^{(l)} \sim p(\epsilon)$ (6)

Auto-Encoding Variational Bayesian algorithm

Algorithm 1 Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings $M = 100$ and $L = 1$ in experiments.

$\theta, \phi \leftarrow$ Initialize parameters

repeat

$\mathbf{X}^M \leftarrow$ Random minibatch of M datapoints (drawn from full dataset)

$\epsilon \leftarrow$ Random samples from noise distribution $p(\epsilon)$

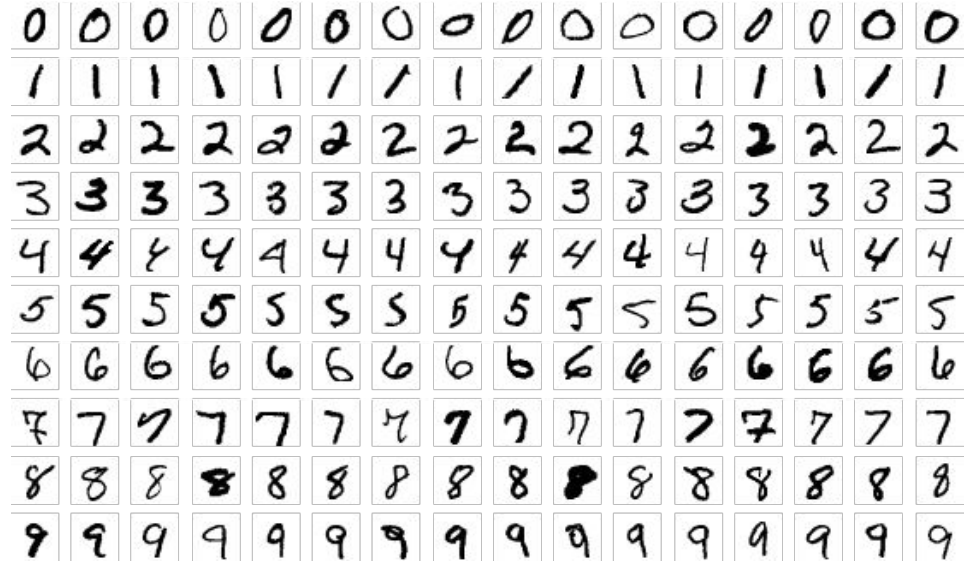
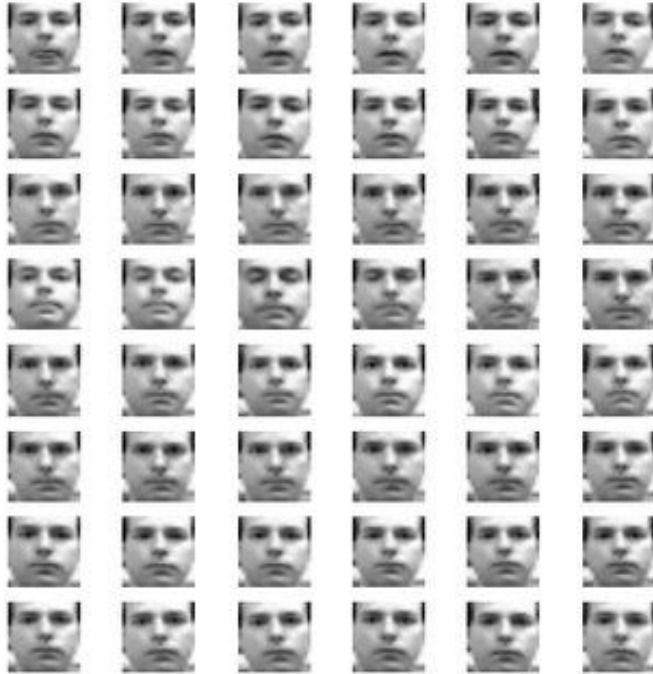
$\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$ (Gradients of minibatch estimator (8))

$\theta, \phi \leftarrow$ Update parameters using gradients \mathbf{g} (e.g. SGD or Adagrad [DHS10])

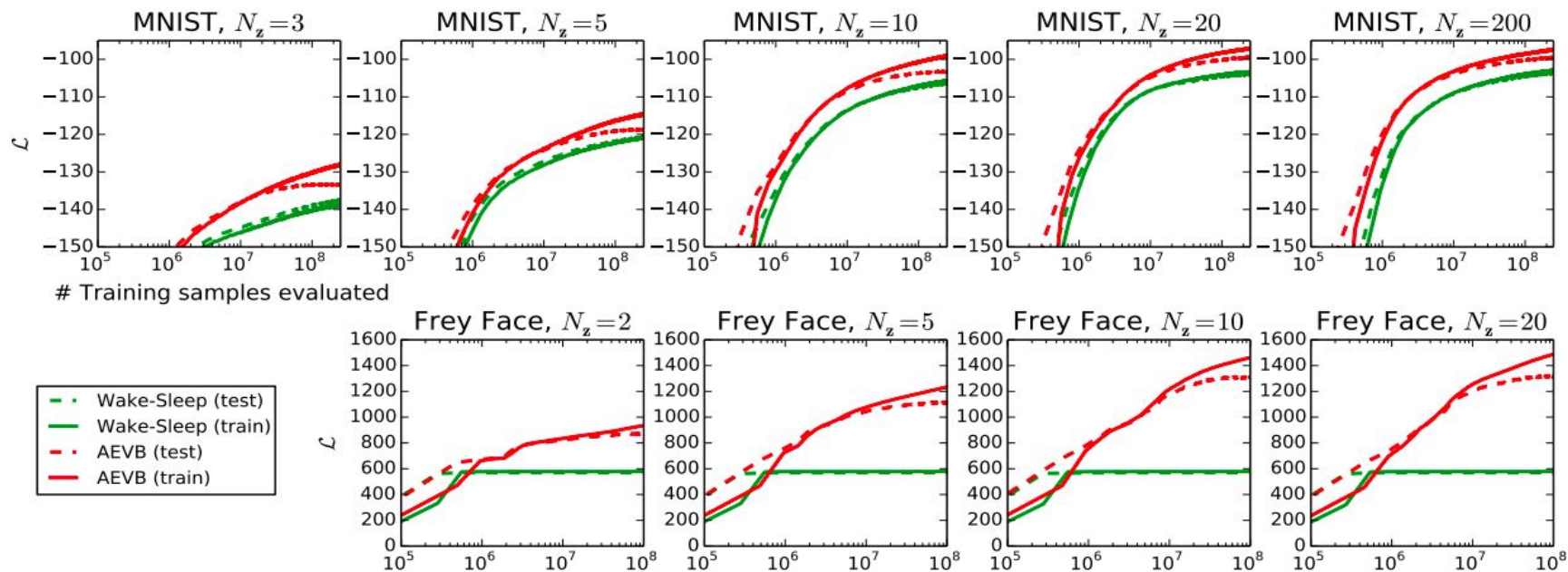
until convergence of parameters (θ, ϕ)

return θ, ϕ

Frey Face and MNIST Datasets



Results for Frey Face and MNIST:



Visualisations of learned data manifold for generative models with two-dimensional latent space, learned with AEVB.



(a) 2-D latent space

(b) 5-D latent space

(c) 10-D latent space

(d) 20-D latent space

Thank you for attention

