Quantum Variational Autoencoder and its applications

Hossein Sadeghi September 16th, 2019 TRIUMF-Helmholtz Workshop



Overview

Quantum Annealing

Physics discovery

Generative modeling with Quantum Annealing

Performance of simulated quantum variational autoencoder

Quantum Annealing

Quantum annealing

- Superconducting circuit as a qubit
- effective Hamiltonian

$$H_p = \sum_i h_i \sigma_i^z + \sum_{ij} J_{ij} \sigma_i^z \sigma_j^z$$
(1)

quantum annealing

$$H = A(s)\sum_{i}\sigma_{i}^{x} + B(s)H_{p}$$

- 2048 qubits, 6k couplers
- Next generation 5k qubits, 40K couplers



Optimization & Sampling

 Many complex problems can be mapped to energy optimization over binary variables

$$E(\mathbf{s}) = \sum_{i} h_i s_i + \sum_{ij} J_{ij} s_i s_j$$

$$\mathbf{s} = \operatorname*{argmin}_{\mathbf{s}} \sum_{i} h_{i} s_{i} + \sum_{ij} J_{ij} s_{i} s_{j}$$

 Sampling. Quantum annealer is a stochastic sampler and it follows "a" Boltzmann distribution.

$$p(\mathbf{s}) = \frac{e^{-\beta E(\mathbf{s})}}{\sum_{\mathbf{s}'} e^{-\beta E(\mathbf{s}')}}$$
$$\mathbf{s} \sim p(\mathbf{s})$$

3/33 Copyright © D-Wave Systems Inc.

Physics discovery

Programmable quantum spin glass simulator



PHYSICS

Phase transitions in a programmable quantum spin glass simulator

R. Harris^{1*}, Y. Sato¹, A. J. Berkley¹, M. Reis¹, F. Altomare¹, M. H. Amin^{1,2}, K. Boothby¹, P. Bunyk¹, C. Deng¹, C. Enderud¹, S. Huang¹, E. Hoskinson¹, M. W. Johnson¹, E. Ladizinsky¹, N. Ladizinsky¹, T. Lanting¹, R. Li¹, T. Medina¹, R. Molavi^{1,3}, R. Neufeld¹, T. Oh¹, I. Pavlov¹, I. Perminov¹, G. Poulin-Lamarre¹, C. Rich¹, A. Smirnov¹, L. Swenson¹, N. Tsai¹, M. Volkmann¹, J. Whittaker¹, J. Yao¹

Material simulation

nature International journal of science

Letter Published: 22 August 2018

Observation of topological phenomena in a programmable lattice of 1,800 qubits

Andrew D. King 🎮, Juan Carrasquilla, […] Mohammad H. Amin

Higgs Boson optimization problem



Letter | Published: 18 October 2017

Solving a Higgs optimization problem with quantum annealing for machine learning

Alex Mott, Joshua Job, Jean-Roch Vlimant, Daniel Lidar & Maria Spiropulu 🏁

Generative modeling with Quantum Annealing

Quantum Boltzmann Machine

 A Boltzmann machine with a transverse field Ising Hamiltonian, with

$$\Gamma = A/B$$

$$H = \Gamma \sum_{i} \sigma_{i}^{x} + \left(\sum_{i} h_{i} \sigma_{i}^{z} + \sum_{i,j} J_{ij} \sigma_{i}^{z} \sigma_{j}^{z} \right)$$
(2)

Objective function: negative log-Likelihood function

$$\mathcal{L} = -\mathbb{E}_{\mathbf{x} \sim D} \log(p(\mathbf{x})) \tag{3}$$

Training:

$$\delta h_{i} = \eta \left(\overline{\langle \sigma_{i}^{z} \rangle_{\mathbf{x}}} - \langle \sigma_{i}^{z} \rangle \right), \qquad (4)$$

$$\delta J_{ij} = \eta \left(\overline{\langle \sigma_{i}^{z} \sigma_{j}^{z} \rangle_{\mathbf{x}}} - \langle \sigma_{i}^{z} \sigma_{j}^{z} \rangle \right). \qquad (5)$$

Quantum Boltzmann Machine

- For a powerful RBM, NM connections are needed. N binary variables = size of binary feature space. M hidden units. For the MNIST dataset with 784 variables and an RBM with 100 hidden units, the number of pairwise couplers is 78400. A bipartite graph with at least the degree 784.
- QBM only works with binary data, or data needs to be binarized.
- A quantum bound (Golden-Thompson inequality) is needed to computed clamped expectations.

Quantum Variational Autoencoders

- A discrete variational autoencoder with a BM or QBM prior
- Quantum annealing is used for sampling from the Boltzmann prior
- The approximate inference is a factorial Bernoulli distribution
- Intractable likelihood evaluation, but a tractable evidence lower bound for optimization

$$\log p_{\theta}(\boldsymbol{x}) \geq \mathbb{E}_{z \sim q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \log \frac{p_{\theta}(\boldsymbol{x}|\boldsymbol{z})p_{\theta}(\boldsymbol{z})}{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})}$$
(6)



Quantum Variational Autoencoders

- With the right encoder and decoder, QVAE can be applied to any type of data
- The number of qubits and the connectivity of the qubits can be controlled.

Back-propagation with discrete variables

$$\nabla_{\theta} \mathbb{E}_{z \sim q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \log \frac{p_{\theta}(\boldsymbol{x}|\boldsymbol{z})p_{\theta}(\boldsymbol{z})}{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})}$$
(7)
= $\mathbb{E}_{z \sim q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \nabla_{\theta} \log \frac{p_{\theta}(\boldsymbol{x}|\boldsymbol{z})p_{\theta}(\boldsymbol{z})}{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})}$ (8)

$$\nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \log \frac{p_{\theta}(\boldsymbol{x}|\boldsymbol{z})p_{\theta}(\boldsymbol{z})}{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})}$$
(9)

Reparameterization trick:

$$\nabla_{\phi} \mathbb{E}_{\rho \sim \mathcal{U}} \log \frac{p(\boldsymbol{x} | \boldsymbol{z}(\rho, \phi)) p(\boldsymbol{z}(\rho, \phi))}{q(\boldsymbol{z}(\rho, \phi) | \boldsymbol{x})}$$
(10)
= $\mathbb{E}_{\rho \sim \mathcal{U}} \nabla_{\phi} \log \frac{p(\boldsymbol{x} | \boldsymbol{z}(\rho, \phi)) p(\boldsymbol{z}(\rho, \phi))}{q(\boldsymbol{z}(\rho, \phi) | \boldsymbol{x})}$ (11)

Back-propagation with discrete variables

The gradient of inverse CDF for discrete variables is undefined.

$$\mathbb{E}_{\rho \sim \mathcal{U}} \nabla_{\phi} \log \frac{p(\mathbf{x} | \mathbf{z}(\rho, \phi)) p(\mathbf{z}(\rho, \phi))}{q(\mathbf{z}(\rho, \phi) | \mathbf{x})}$$
(12)

Continuous relaxation $z \to \zeta.$ Note: biased low-variance estimate of the gradient.

$$\mathbb{E}_{\rho \sim \mathcal{U}} \nabla_{\phi} \log \frac{p(\mathbf{x} | \zeta(\rho, \phi)) p(\zeta(\rho, \phi))}{q(\zeta(\rho, \phi) | \mathbf{x})}$$
(13)

$$\zeta(\rho,\phi) = \sigma(\frac{l_{\phi}+\rho}{\tau}), \ \rho \sim \mathcal{U}$$
(14)

Performance of simulated quantum variational autoencoder

Generative performance of QVAE+QBM

29026982931743 86778017867516 04507402190001 74036041492233 01084210168470 10028492121441 18107511962550 28179730850603 47144181763300 15312049409014 01968/19814154 41319733220903 48823029048433 47622407825662 76753432035617 29517863014053 06723076468568 78361836070660 77998146981751 69727244409626 33978280310344 14969776975913 36458558799184 25338638362117 96983691171990 79022892770719 22688677288425 82366388920202 (b) Generated MNIST: QVAE, (d) Generated MNIST: QVAE, QBM_{64×64}, $\Gamma = 0$, QBM_{64×64}, $\Gamma = 1$.

Quantitative performance

MNIST (static binarization)				
		ELBO	Q-ELBO	
QVAE: Г	$= 0 \operatorname{RBM}_{16 \times 16}$	-109.3 ± 0.2	-109.3 ± 0.2	
Г	$= 1 \text{ QBM}_{16 \times 16}$	-110.5	-120.6	
Г	= 2	-115.3	-135.8	
QVAE: Γ	$= 0 \operatorname{RBM}_{32 \times 32}$	-101.8	-101.8	
Γ	$= 1 \text{ QBM}_{32 \times 32}$	-103.6	-117.9	
Γ	= 2	-112.1	-139.7	
QVAE: Γ	$= 0 \operatorname{RBM}_{64 \times 64}$	-105.7	-105.7	
Γ	$= 1 \text{ QBM}_{64 \times 64}$	-108.7	-133.9	
Г	= 2	-120.0	-165.2	

Generative performance of QVAE+RBM

78852966884411 78150719489247 19263904760136 59927992490497 13969140589831 34274811510715 1130475424364 72996999494973 364660112#63056 1958799**4**997**3**99 22432677758448 42697049718374 71095181897119 79348474440597 91668386738037 04786447985883 78197740019496 9391991/1917392 91706981974086 08077043412142 239536/2102002 19419511998494 67440841385936 87413734947949 18647893914834 74839790490394 73396967994263 29319194264682 (a) Generated MNIST: DVAE, (c) Generated MNIST: DVAE, $RBM_{32\times32}$ $RBM_{256\times 256}$

Quantitative performance

MNIST (static binarization)				
		ELBO	$\mathbf{L}\mathbf{L}$	
DVAE	$\mathrm{RBM}_{32 \times 32}$	-99.3 ± 0.2	-90.8 ± 0.2	
	$\mathrm{RBM}_{64 \times 64}$	-92.4	-85.5	
	$\mathrm{RBM}_{128\times 128}$	-90.4	-84.7	
	$\mathrm{RBM}_{256\times 256}$	-89.2	-83.5	
VIMCO [88]			-91.9	
NVIL [75]			-93.5	
CONCRETE [82]			-85.7	
GS [76]		-101.5		
RWS [89]			-88.9	
REBAR [90]		-98.8		

PixelQVAE

- Combination of a QVAE with PixelCNN decoder
- Autoregressive decoder, latent variable model, variational inference

$$\log p(\boldsymbol{x}|\boldsymbol{z}) = \sum_{i} \log p(x_i|\boldsymbol{x}_{< i}, \boldsymbol{z})$$
(15)

$$\log p(\mathbf{x}) \ge \mathbb{E}_{z \sim q(\mathbf{z}|\mathbf{x})} \log \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q(\mathbf{z}|\mathbf{x})}$$
(16)

Results

	MNIST (Static)		MNIST (Dynamic)	
	NLL	KL	NLL	KL
PixelVAE++ Gaussian	78.66	6.86	78.01	4.2
PixelVAE++ RBM	78.65	7.62	78.00	5.05
VLAE	79.03		78.53	
	OMNIGLOT		Caltech 101	
	NLL	KL	NLL	KL
PixelVAE++ Gaussian	88.65	1.63	79.52	4.00
PixelVAE++ RBM	88.29	2.56	77.46	6.85
VLAE	89.83		77.36	
bpd	CIFAR10			
	NLL	KL		
PixelVAE++ Gaussian	2.92	0.005		
PixelVAE++ RBM	2.90	0.016		
VLAE	2.95			
PixelCNN++	2.92			
PixelSNAIL	2.85			

Summary

- Powerful latent variables models can be hybridized with Quantum (assisted) Boltzmann Machines.
- The range of applications extends to that of all VAE models
- Quantum annealing can circumvent the issue of sampling from Boltzmann distributions and replace expensive MCMC methods.
- Currently, quantum annealing can achieve similar performance
- As the decoder gets more powerful, the role of latent variables fades away
- Sampling bias can prevent gaining any advantage from a QVAE.

Thank you!

Quantum Boltzmann Machine

- A hybrid ML algorithm that makes use of deep neural networks for encoding and decoding latent variables
- The prior distribution is a classical/quantum Boltzmann machine

$$E_{\mathbf{z}} = -\sum_{a} b_a z_a - \sum_{a,b} w_{ab} z_a z_b.$$
(17)

$$P_{\mathbf{v}} = Z^{-1} \sum_{\mathbf{h}} e^{-E_{\mathbf{z}}}, \qquad Z = \sum_{\mathbf{z}} e^{-E_{\mathbf{z}}},$$
(18)

$$\mathcal{L} = -\sum_{\mathbf{v}} P_{\mathbf{v}}^{\text{data}} \log P_{\mathbf{v}}, \tag{19}$$

$$\mathcal{L} = -\sum_{\mathbf{v}} P_{\mathbf{v}}^{\text{data}} \log \frac{\sum_{\mathbf{h}} e^{-E_{\mathbf{z}}}}{\sum_{\mathbf{z}'} e^{-E_{\mathbf{z}'}}}.$$
 (20)

$$\delta\theta = -\eta \partial_{\theta} \mathcal{L}, \tag{21}$$

20/33 Copyright © D-Wave Systems Inc.

D:Wave

$$\partial_{\theta} \mathcal{L} = \sum_{\mathbf{v}} P_{\mathbf{v}}^{\text{data}} \frac{\sum_{\mathbf{h}} \partial_{\theta} E_{\mathbf{z}} e^{-E_{\mathbf{z}}}}{\sum_{\mathbf{h}} e^{-E_{\mathbf{z}}}} - \frac{\sum_{\mathbf{z}} \partial_{\theta} E_{\mathbf{z}} e^{-E_{\mathbf{z}}}}{\sum_{\mathbf{z}} e^{-E_{\mathbf{z}}}}$$
$$= \overline{\langle \partial_{\theta} E_{\mathbf{z}} \rangle_{\mathbf{v}}} - \langle \partial_{\theta} E_{\mathbf{z}} \rangle, \qquad (22)$$

$$\delta b_a = \eta \left(\overline{\langle z_a \rangle_{\mathbf{v}}} - \langle z_a \rangle \right), \qquad (23)$$

$$\delta w_{ab} = \eta \left(\overline{\langle z_a z_b \rangle_{\mathbf{v}}} - \langle z_a z_b \rangle \right).$$
 (24)

$$H = -\sum_{a} b_a \sigma_a^z - \sum_{a,b} w_{ab} \sigma_a^z \sigma_b^z.$$
 (25)

$$\sigma_a^z \equiv \overbrace{I \otimes \dots \otimes I}^{a-1} \otimes \sigma_z \otimes \overbrace{I \otimes \dots \otimes I}^{N-a}$$
(26)

$$I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \qquad \sigma_z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}.$$
(27)

$$\rho = Z^{-1} e^{-H}.$$
(28)

$$P_{\mathbf{v}} = \mathsf{Tr}[\Lambda_{\mathbf{v}}\rho],\tag{29}$$

$$\Lambda_{\mathbf{v}} = \mathbf{v}\mathbf{v}\otimes\mathcal{I}_{\mathbf{h}},\tag{30}$$

$$\mathbf{v}\mathbf{v} \equiv \prod_{\nu} \left(\frac{1 + \mathbf{v}_{\nu} \sigma_{\nu}^{z}}{2} \right)$$
(31)

$$\sigma_a^x \equiv \overbrace{I \otimes \ldots \otimes I}^{a-1} \otimes \sigma_x \otimes \overbrace{I \otimes \ldots \otimes I}^{N-a}, \qquad \sigma_x = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix},$$

$$\mathcal{L} = -\sum_{\mathbf{v}} P_{\mathbf{v}}^{\text{data}} \log \frac{\text{Tr}[\Lambda_{\mathbf{v}} e^{-H}]}{\text{Tr}[e^{-H}]}.$$
(32)

$$\partial_{\theta} \mathcal{L} = \sum_{\mathbf{v}} P_{\mathbf{v}}^{\text{data}} \left(\frac{\text{Tr}[\Lambda_{\mathbf{v}} \partial_{\theta} e^{-H}]}{\text{Tr}[\Lambda_{\mathbf{v}} e^{-H}]} - \frac{\text{Tr}[\partial_{\theta} e^{-H}]}{\text{Tr}[e^{-H}]} \right).$$
(33)

$$\partial_{\theta} e^{-H} = \sum_{m=1}^{n} e^{-m\delta\tau H} \left(-\partial_{\theta} H \delta\tau\right) e^{-(n-m)\delta\tau H}.$$
(34)

$$\partial_{\theta} e^{-H} = -\int_0^1 d\tau e^{-\tau H} \partial_{\theta} H e^{(\tau-1)H}.$$
 (35)

$$\operatorname{Tr}[\partial_{\theta}e^{-H}] = -\operatorname{Tr}[\partial_{\theta}He^{-H}], \qquad (36)$$

$$\frac{\mathrm{Tr}[\partial_{\theta}e^{-H}]}{\mathrm{Tr}[e^{-H}]} = -\langle \partial_{\theta}H \rangle, \tag{37}$$

$$\frac{\mathrm{Tr}[\Lambda_{\mathbf{v}}\partial_{\theta}e^{-H}]}{\mathrm{Tr}[\Lambda_{\mathbf{v}}e^{-H}]} = -\int_{0}^{1} dt \frac{\mathrm{Tr}[\Lambda_{\mathbf{v}}e^{-tH}\partial_{\theta}He^{-(1-t)H}]}{\mathrm{Tr}[\Lambda_{\mathbf{v}}e^{-H}]}, \qquad (38)$$

$$\operatorname{Tr}[e^{A}e^{B}] \ge \operatorname{Tr}[e^{A+B}], \tag{39}$$

$$P_{\mathbf{v}} = \frac{\operatorname{Tr}[e^{-H}e^{\ln\Lambda_{\mathbf{v}}}]}{\operatorname{Tr}[e^{-H}]} \ge \frac{\operatorname{Tr}[e^{-H+\ln\Lambda_{\mathbf{v}}}]}{\operatorname{Tr}[e^{-H}]}.$$
(40)

$$H_{\mathbf{v}} = H - \ln \Lambda_{\mathbf{v}},\tag{41}$$

$$P_{\mathbf{v}} \ge \frac{\mathrm{Tr}[e^{-H_{\mathbf{v}}}]}{\mathrm{Tr}[e^{-H}]}.$$
(42)

$$H_{\mathbf{v}} \equiv H(\sigma_{\nu}^{\mathcal{X}} = 0, \sigma_{\nu}^{\mathcal{Z}} = \mathbf{v}_{\nu}).$$
(43)

$$\mathcal{L} \leq \tilde{\mathcal{L}} \equiv -\sum_{\mathbf{v}} P_{\mathbf{v}}^{\text{data}} \log \frac{\text{Tr}[e^{-H_{\mathbf{v}}}]}{\text{Tr}[e^{-H}]}.$$
(44)

$$\partial_{\theta} \tilde{\mathcal{L}} = \sum_{\mathbf{v}} P_{\mathbf{v}}^{\text{data}} \left(\frac{\text{Tr}[e^{-H_{\mathbf{v}}} \partial_{\theta} H_{\mathbf{v}}]}{\text{Tr}[e^{-H_{\mathbf{v}}}]} - \frac{\text{Tr}[e^{-H} \partial_{\theta} H]}{\text{Tr}[e^{-H}]} \right),$$

$$= \left(\overline{\langle \partial_{\theta} H_{\mathbf{v}} \rangle_{\mathbf{v}}} - \langle \partial_{\theta} H \rangle \right), \qquad (45)$$

where

$$\overline{\langle ... \rangle_{\mathbf{v}}} = \sum_{\mathbf{v}} P_{\mathbf{v}}^{\text{data}} \langle ... \rangle_{\mathbf{v}} = \sum_{\mathbf{v}} P_{\mathbf{v}}^{\text{data}} \frac{\text{Tr}e^{-H_{\mathbf{v}}}...}{\text{Tr}e^{-H_{\mathbf{v}}}}.$$
(46)

Taking heta to be b_a , w_{ab} , and using $\delta heta = -\eta \partial_{ heta} \tilde{\mathcal{L}}$, we obtain

$$\delta b_a = \eta \left(\overline{\langle \sigma_a^z \rangle_{\mathbf{v}}} - \langle \sigma_a^z \rangle \right), \tag{47}$$

$$\delta w_{ab} = \eta \left(\overline{\langle \sigma_a^z \sigma_b^z \rangle_{\mathbf{v}}} - \langle \sigma_a^z \sigma_b^z \rangle \right).$$
(48)

$$\delta\Gamma_a = \eta \left(\overline{\langle \sigma_a^x \rangle_{\mathbf{v}}} - \langle \sigma_a^x \rangle \right). \tag{49}$$

$$H_{\mathbf{v}} = -\sum_{i} \left(\Gamma_{i} \sigma_{i}^{x} + b_{i}^{\text{eff}}(\mathbf{v}) \sigma_{i}^{z} \right),$$
(50)

where $b_i^{\text{eff}}(\mathbf{v}) = b_i + \sum_{\nu} w_{i\nu} \mathbf{v}_{\nu}$. Expectations $\langle \sigma_i^z \rangle_{\mathbf{v}}$ entering (47) can be computed exactly:

$$\langle \sigma_i^z \rangle_{\mathbf{v}} = \frac{b_i^{\text{eff}}}{D_i} \tanh D_i,$$
 (51)

where $D_i = \sqrt{\Gamma_i^2 + (b_i^{\rm eff})^2}$. Notice that (51) reduces to the classical RBM expression,

$$\langle \sigma_i^z \rangle_{\mathbf{v}} = \tanh b_i^{\text{eff}},$$
 (52)

Results

	MNIST (Static)		MNIST (Dynamic)	
	NLL	KL	NLL	KL
PixelVAE++ Gaussian	78.66	6.86	78.01	4.2
PixelVAE++ RBM	78.65	7.62	78.00	5.05
VLAE	79.03		78.53	
	OMNIGLOT		Caltech 101	
	NLL	KL	NLL	KL
PixelVAE++ Gaussian	88.65	1.63	79.52	4.00
PixelVAE++ RBM	88.29	2.56	77.46	6.85
VLAE	89.83		77.36	
bpd	CIFAR10			
	NLL	KL		
PixelVAE++ Gaussian	2.92	0.005		
PixelVAE++ RBM	2.90	0.016		
VLAE	2.95			
PixelCNN++	2.92			
PixelSNAIL	2.85			

Mutual information

- VAEs are notorious for deactivating latent variables
- various methods are used to overcome this issue, but when decoder gets stronger (PixelVAE) things get worse
- This issue can be expressed in terms of mutual information between the observations and latent variables. (quantified in terms of KL)
- Training PixelVAE naively will result in a redundant VAE
- Scheduled sampling, noisy auto-regressive channel, warm-up, ...
- Using discrete variables seems to help

Image generation



Graph visualization



Figure: (a) A VAE model, (b) An autoregressive model, (c) A VAE with an autoregressive decoder, and (d) Inference model.

PixelVAE++



(a) PixelVAE++ architecture

D:Wave

PixelCNN (an Autoregressive model)

- Factorizes the joint distribution of the features to a product of conditionals
- Tractable likelihood. Fast training. High capacity model
- ► Slow sample generation. There are ways to improve.
- Conditional probabilities are parameterized by masked convolutional neural networks with shared parameters
- Useful when there are local temporal or spatial correlations

$$\log p(\mathbf{x}) = \sum_{i} \log p(x_i | \mathbf{x}_{< i})$$
(53)

