

Approximate Inference with Amortised MCMC

Yingzhen Li

MLG, University of Cambridge

Joint work with Rich Turner (Cambridge) and Qiang Liu (Dartmouth) arXiv preprint 1702.08343



Recap: from Variational EM to VAE



(For simplicity we omit the model parameter θ but it would be trained by approx. MLE)

Recap: from Variational EM to VAE



Amortised inference: memory efficient & no need to run VI optimisation in test time!



For every x_n , need to simulate MCMC for T >> 0 steps (slow!)



Need to store all samples from the previous iteration, memory cost O(NKD). For a new datapoint, still need to run MCMC with T >> 0 starting from p_0 (slow!)



This method essentially cares about q_T only, so no need for $q(z|x) \approx p(z|x)$. In test time still need to run MCMC to obtain samples from q_T (slow!)



Distillation happens at the same time during training (thus also improving q_T), and now $q(z|x) \approx p(z|x)$ – no need to run MCMC in test time!

Understanding "distillation during training"



(think about *projected* gradient descent)

We want q to be implicit!

(i.e. can sample from q but cannot evaluate density)

idea: match samples $\{ \pmb{z}_0^k \} \sim q$ to samples $\{ \pmb{z}_T^k \} \sim q_T !$

- We tested the original GAN idea 1
- In general, any GAN-like technique is applicable!

¹Goodfellow et al. Generative Adversarial Networks. NIPS 2014.

Generative model: a small convolutional decoder

- VAE training:
 - Gaussian encoder: a symmetric flip of the decoder
- Amortised MCMC (AMC) training:
 - CNN-G encoder: $z = MLP([CNN(x), \epsilon]), \epsilon \sim \mathcal{N}(\mathbf{0}, I)$
 - CNN-B encoder: $z = MLP([CNN(x) \odot \epsilon]), \epsilon \sim Bern(0.5)$
 - MCMC: Langevin Dynamics w/out rejection

0/23+56789	2 0 0 0 4 7 5 5 4
0 23456789	6 7 5 7 8 9 5 9 2 9
0 23456789	7 3 4 0 5 9 4 1 9
0 23456789	7 4 8 4 5 4 1 7 5
0 23456789	7 4 9 3 4 5 4 1 9
0 23456789	7 4 7 3 9 2 7 2 9
0/23456789	9 3 4 0 5 1 9 2 7 2 9
0/23456789	7 4 7 6 4 5 4 8 0 2
0/23456789	3 5 8 6 4 8 9 2 4 9
0/23456789	7 4 7 6 9 1 9 6 1 0
0 23456789	Gaussian encoder + VAE
7 G S <i>F</i> S 7 G R S S 7 G S S 7 G R S S 7 G S S S S 2 / B S I S G S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S S	0 0 47 37 37 37 41 3 9 0 5 3 4 2 3 9 3 7 1 9 1 2 3 4 2 3 5 7 5 3 6 7 4 4 8 1 7 4 9 6 9 2 4 8 7 4 2 5 3 0 3 2 7 9 6 8 1 4 5 5 3 0 1 4 9 6 7 8 6 9 3 6 9 8 4 4 3 7 4 // 2 1 4 5 9 6 6 9 2 6 9 7 7 2 CNN-B + AMC

Table 1: Average Test LL and effective sample size (ESS). η as the stepsize for Langevin dynamics.

Encoder	Method	IS-LL	IS-ESS	HAIS-LL	HAIS-ESS
Gaussian	VAE	-81.31	104.11	-80.64	91.59
	MCMC-VI, $T=$ 5, $\eta=$ 0.2	-90.06	110.58	-89.79	85.63
	AMC, $T=$ 5, $\eta=$ 0.2	-90.71	49.02	-89.64	87.93
CNN-G	AMC, $T=$ 5, $\eta=$ 0.2	-90.84	31.60	-89.35	87.49
	AMC, $T=$ 50, $\eta=$ 0.02	-83.30	6.84	-78.23	77.78
	AVB	-94.97	11.30	-85.92	57.21
CNN-B	AMC, $T=$ 5, $\eta=$ 0.2	-90.75	34.17	-89.42	88.10
	AMC, $T=50$, $\eta=0.02$	-83.62	8.88	-80.03	80.71
	AVB	-89.47	8.98	-82.66	76.90
N/A	persistent MCMC, $T=50,~\eta=0.02$	-84.43	9.14	-78.88	77.29

• HAIS seems to be more reliable (K = 100), compared to importance sampling (IS, K = 5000)

• The best case (CNN-G) is better than persistent MCMC

Missing data imputation

- Given an image $\mathbf{x} = [\mathbf{x}_o, \mathbf{x}_m]$ with missing values, repeat the following for T steps:
 - sample $z \sim q(z|x_o, x_m)$
 - sample $x^* \sim p(x|z)$ and set $x_m \leftarrow x_m^*$

Table 2: Label entropy on nearest neighbours. The l_1 distance is divided by the number of pixels.

Dataset	VAE	CNN-G	CNN-B
Entropy	$0.411{\pm}0.0389$	$0.701{\pm}0.0476$	$0.933{\pm}0.0491$
<i>l</i> ₁ -norm	$0.061{\pm}0.0002$	$0.059 {\pm} 0.0001$	$0.064{\pm}0.0002$



- A lot more to be done!
 - what's the best combo of q, MCMC algorithm, and the distillation rule? (Langevin dynamics & original GAN are inefficient, now trying HMC & WGAN)
- Reuse intermediate samples from q_1 to q_{T-1} ?
- Discrete distributions?
- Other ideas to mix MCMC, VI and implicit distributions?
 - e.g. see "reparameterised MCMC" by Michalis Titsias

Come and find me at the posters! (we have more distillation rules & results) (also see later spotlight on gradient estimators)