Approximate Inference and Generative Models

Danilo J. Rezende

Google DeepMind

Hammers & Nails - Machine Learning & HEP July 19-28, 2017 | Weizmann Institute of Science, Israel

Making DNNs a little bit more Bayesian



Weight Uncertainty in Neural Networks https://arxiv.org/abs/1505.05424

Why Generative Models?

- Compact mechanism to express causal relations and dependencies
- Make predictions (with uncertainty)
- Plan sequences of decisions (planning as inference)
- Experiment design (planning + expected information gain)
- Data compression
- Hypothesis testing
- Handle missing data

Definitions

Observed variables

Unobserved or latent variables

$$\begin{aligned} &\{x_i\} \quad x_i \in \Omega \text{ (e.g. } \mathbb{R}^{d_x}, [0,1]^{d_x}, \{0,1\}^{d_x}) \\ &\{z_i\} \quad z_i \in \Lambda \text{ (e.g. } \mathbb{R}^{d_z}) \\ &\theta \quad \in \mathbb{R}^{d_\theta} \end{aligned}$$

$$\begin{array}{ll} p &: \Omega \otimes \Lambda \otimes \mathbb{R}^{d_{\theta}} \to \mathbb{R}^{+} \quad s.t. \int_{x \in \Omega, z \in \Lambda, \theta \in \mathbb{R}^{m}} dx \, dz d\theta p(x, z, \theta) = 1 \\ p &\in C^{2} \end{array}$$

$$p(x, z, \theta) = \frac{1}{\mathcal{Z}} e^{-U(x, z, \theta)}$$
$$\mathbb{E}_{x \sim p}[f(x)] = \int dx p(x) f(x)$$

Graphical Models



$p(x, z, \theta) = \rho(\theta) \prod_{i=1}^{N} p(x_i | z_i, \theta) \pi(z_i)$

Graphical Models



Graphical Models: a concrete example



$$p(x, z, \theta) = \rho(\theta) \prod_{i=1}^{N} p(x_i | z_i, \theta) \pi(z_i)$$
$$\pi(z) = \mathcal{N}(0, \mathbb{I}_{d_z})$$
$$\rho(\theta) = \mathcal{N}(0, \kappa^2 \mathbb{I}_{d_\theta})$$
$$p(x | z, \theta) = \mathcal{N}(\theta_0 + \theta_1 z, \exp(\theta_2))$$

 $\theta = \{\theta_0 \in \mathbb{R}^{d_x}, \theta_1 \in \mathbb{R}^{d_x \times d_z}, \theta_2 \in \mathbb{R}^{d_x}\}$

Graphical Models + Computational Graphs (aka NNets)



$$\begin{aligned} \pi(z) &= \mathcal{N}(0, \mathbb{I}_{d_z}) \\ \rho(\theta) &= \mathcal{N}(0, \kappa^2 \mathbb{I}_{d_\theta}) \\ p(x|z, \theta) &= \mathcal{N}(\theta_0 + \theta_1 z, \exp(\theta_2)) \\ \pi(z) &= \mathcal{N}(0, \mathbb{I}_{d_z}) \\ \rho(\theta) &= \mathcal{N}(0, \kappa^2 \mathbb{I}_{d_\theta}) \\ h_1 &= \theta_0 + \theta_1 z \\ h_2 &= \exp(\theta_2) \\ p(x|z, \theta) &= \mathcal{N}(h_1, h_2) \end{aligned}$$

Graphical Models + Computational Graphs (alternative notation)

$$\begin{array}{c|c} & & & & & \\ \hline \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \hline \\$$

What do we want to know about $p(x, z, \theta)$?

Goal	Quantity of Interest
Prediction	$p(x_{(t+1),\dots,\infty} x_{-\infty,\dots,t})$
Planning	$J = \mathbb{E}_p\left[\int_0^\infty dt C(x_t) \middle x_0, u\right]$
Parameter estimation	$p(heta x_{0,,N})$
Experiment Design	$EIG = D[p(f(x_{t,\dots,\infty}) u); p(f(x_{-\infty,\dots,t}))]$
Hypothesis testing	$\frac{p(f(x_{-\infty,,t}) H_0)}{p(f(x_{-\infty,,t}) H_1)}$

Density estimation through conditioning

We observe N samples
$$\,\{x_{i=1,...,N}\}\,$$
 and we want to know $\,\,x_{N+1}\!\sim\!?\,$



$$\mathcal{Z} = p(x_{i=1,\dots,N})$$

Density estimation through conditioning

Let's approximate some integrals...

$$p(x_{N+1}|x_{1,\dots,N}) = \frac{1}{\mathcal{Z}} \int d\theta dz \rho(\theta) p(x_{N+1}, z_{N+1}|\theta) \prod_{i=1}^{N} p(x_i|z_i, \theta) \pi(z_i)$$
$$= \frac{1}{\mathcal{Z}} \int d\theta dz_{N+1} p(x_{N+1}, z_{N+1}|\theta) \rho(\theta) \int dz \prod_{i=1}^{N} p(x_i, z_i|\theta)$$

Density estimation through conditioning

Let's approximate some integrals...





Approximate Inference: importance sampling

Let's approximate some integrals...

$$\{z_k\} \sim q(z)$$

$$M(x,\theta) \approx -\ln \sum_{k=1}^{K} e^{-\mathcal{F}(x,z_k,\theta)} + \ln K$$

Approximate Inference: Laplace/saddle point



$$u(x,z,\theta) \approx u(x,\mu,\theta) + J^T(z-\mu) + \frac{1}{2}(z-\mu)^T H(z-\mu)$$

$$\int dz e^{-u(x,z,\theta)} \approx e^{-u(x,\mu,\theta)} \int dz e^{-J^T(z-\mu) - \frac{1}{2}(z-\mu)^T H(z-\mu)}$$
$$\approx e^{-u(x,\mu,\theta) + \frac{1}{2}J^T H J} \det(2\pi H)^{\frac{1}{2}}$$

$$M(x,\theta) ~\approx u(x,\mu,\theta) - \frac{1}{2}J^T H J - \frac{1}{2} \ln \det(2\pi H)$$

Approximate Inference: Perturbative expansion

$$e^{-M(x,\theta)} = \mathbb{E}_{z \sim q}[e^{-\mathcal{F}(x,z,\theta)}]$$



$$e^{-M(x,\theta)} = \sum_{k=0}^{\infty} \frac{(-1)^k}{k!} \mathbb{E}_{z \sim q} [\mathcal{F}(x,z,\theta)^k]$$

Approximate Inference: Variational approx

$$\begin{split} M(x,\theta) &= -\ln \int dz q(z) e^{-\mathcal{F}(x,z,\theta)} \\ \xrightarrow{\text{Exp is}} &\leqslant \int dz q(z) \mathcal{F}(x,z,\theta) \qquad \forall q \neq 0 \end{split}$$

$$M(x,\theta) \leqslant \operatorname{argmin}_{q} \int dz q(z) \mathcal{F}(x,z,\theta) \quad s.t. \int dz q(z) = 1$$

Approximate Inference: Relation to information theory

 $\mathbb{E}_{z \sim q}[\mathcal{F}(x,$

KLD(q||p) is the number of bits necessary to transmit the distribution q to a receiver that already knowns p

$$(D - \ln \pi(z_i))$$

 $D[q_i | | \pi]]$ gularizer



Approximate Inference: Variational mean-field^e approx

ΛT

$$q(z) = \prod_{i=1}^{N} q_i(z_i) \approx \prod_{i=1}^{N} \prod_{j=1}^{a_z} q_{ij}(z_{ij})$$
$$M(x,\theta) \leq \operatorname{argmin}_{q} \int dz q(z) \mathcal{F}(x,z,\theta)$$

ΛT

1

Fixed-point equations $q_{ij}^{\star}(z_{ij}) \propto e^{-\mathbb{E}_{z_{/ij} \sim q^{\star}}[\mathcal{F}(x,z,\theta)]}$

Variational approx: Amortized inference

N

i=1

q(z)

Requires solving optimization for every new data-point x.

Solution:

 $C(x, z, \theta)$

Parameters phi are shared across all data-points => consistency, fast convergence

$$M(x,\theta) \leqslant \underset{q}{\operatorname{argmin}} \int dz q(z) \mathcal{F}(x,z,\theta) \Rightarrow M(x,\theta) \leqslant \underset{\phi}{\operatorname{argmin}} \int dz \prod_{i} q_{\phi}(z_{i}|x_{i}) \mathcal{F}(x,z,\theta)$$

Variational approx: Amortization & reparametrization

 $\nabla_{\phi} \mathbb{E}_{z_i \sim q_\phi} [\mathcal{F}(x_i, z_i, \theta)] = \nabla_{\phi} \mathbb{E}_{\xi_i \sim \lambda} [\mathcal{F}(x_i, g(\xi_i, \phi), \theta)] = \mathbb{E}_{\xi_i \sim \lambda} [\nabla_{\phi} \mathcal{F}(x_i, g(\xi_i, \phi), \theta)]$

Variational approx: Amortization & reparametrization. Gaussian example:

$$q(z|x) = \mathcal{N}(z|\mu_{\phi}(x), \Sigma_{\phi}(x) = C_{\phi}(x)^{T}C_{\phi}(x))$$

then

$\mathbb{E}_{z \sim q}[f(z)] = \mathbb{E}_{\xi \sim \mathcal{N}(0,\mathbb{I})}[f(\mu_{\phi}(x) + C_{\phi}(x)^{T}\xi)]$

Stochastic Backpropagation and Approximate Inference in Deep Generative Models <u>http://arxiv.org/abs/1401.4082v3</u> Auto-Encoding Variational Bayes <u>https://arxiv.org/abs/1312.6114</u>

Variational approx: Amortization & reparametrization. Gaussian gradients example:

$$\mathbb{E}_{z \sim q}[f(z)] = \mathbb{E}_{\xi \sim \mathcal{N}(0,\mathbb{I})}[f(\mu_{\phi}(x) + C_{\phi}(x)^{T}\xi)]$$

Only using first order information

$$\nabla_{\phi} \mathbb{E}_{\xi \sim \mathcal{N}(0,\mathbb{I})} [f(\mu_{\phi}(x) + C_{\phi}(x)^{T} \xi)] = \mathbb{E}_{\xi \sim \mathcal{N}(0,\mathbb{I})} [J^{T} \{\nabla_{\phi} \mu_{\phi}(x) + \nabla_{\phi} C_{\phi}(x)^{T} \xi\}]$$

Using second order information

$$\mathbb{E}_{z \sim q}[J^T \nabla_{\phi} \mu_{\phi}(x) + \operatorname{Tr}[HC_{\phi}(x) \nabla_{\phi} C_{\phi}(x)]]$$

Stochastic Backpropagation and Approximate Inference in Deep Generative Models <u>http://arxiv.org/abs/1401.4082v3</u> Auto-Encoding Variational Bayes <u>https://arxiv.org/abs/1312.6114</u>

Deep Latent Generative Models / VAES

- Describes a causal process by which data **y** is generated.
 - Layers of stochastic variables **z**

 Z_1

h,

h,

V

n = 1, ..., N

W

- Conditional probabilities using no on deep neural networks.
- Statistical problem of density esti

Latent Variables (Stochastic layers)

$$egin{aligned} \mathbf{z}_l &\sim \mathcal{N}(\mathbf{z}_l | f_l(\mathbf{z}_{l+1}), \mathbf{\Sigma}_l) \ f_l(\mathbf{z}) &= \sigma(\mathbf{W}h(\mathbf{z}) + \mathbf{b}) \end{aligned}$$

Deterministic layers

$$h_i(\mathbf{x}) = \sigma(\mathbf{A}\mathbf{x} + \mathbf{c})$$

Observation Model

 $\eta = Wh_1 + b$ $y \sim Expon(y|\eta)$

Can also use non-exponential family.

Two processes to understand



1. Generate new data

- Sample new data given underlying causes.

2. Explain observed data

- Recognition or inference of observed data to obtain the posterior distribution.

Generative model and posterior representation are connected by variational inference.



Stochastic Backpropagation and Approximate Inference in Deep Generative Models <u>http://arxiv.org/abs/1401.4082v3</u> Auto-Encoding Variational Bayes <u>https://arxiv.org/abs/1312.6114</u>

Demo



Improving Inference

Desired properties of an inference model

Flexible:

• Can increase complexity on demand

Fast at training time:

- Cheap to sample from
- Unbiased approximation to entropy derivatives known

Fast at test time:

• Should not rely on the generative model's likelihoods or derivatives at test time

Richer Families of Posteriors

Two high-level goals:
Build richer approximate posterior distributions.
Maintain computational efficiency and scalability



Same as the problem of specifying a model of the data itself

Richer Families of Posteriors





Richer Families of Posteriors



Amortized inference with Normalizing Flows



Example: Planar Flow in \mathbb{R}^2



Variational Inference with Normalizing Flows https://arxiv.org/abs/1505.05770

Normalizing Flows on non-Euclidean Manifolds

Normalizing Flows on non-Euclidean Manifolds

Normalizing Flows on Riemannian Manifolds https://arxiv.org/abs/1611.02304

Improved Normalizing Flows

(a) Forward propagation

(b) Inverse propagation

$[\mathbf{m}_t, \mathbf{s}_t]$	$\leftarrow \texttt{AutoregressiveNN}$	[t]	$(\mathbf{z}_t,$	h ;θ)
--------------------------------	----------------------------------------	-----	------------------	-------------	---

$$\sigma_t = \text{sigmoid}(\mathbf{s}_t)$$
$$\mathbf{z}_t = \sigma_t \odot \mathbf{z}_{t-1} + (1 - \sigma_t) \odot \mathbf{m}_t$$

Dataset	PixelRNN [46]	Real NVP	Conv DRAW [22]	IAF-VAE [34]
CIFAR-10	3.00	3.49	< 3.59	< 3.28
Imagenet (32×32)	3.86 (3.83)	4.28 (4.26)	< 4.40 (4.35)	
Imagenet (64×64)	3.63 (3.57)	3.98 (3.75)	< 4.10 (4.04)	
LSUN (bedroom)		2.72 (2.70)		
LSUN (tower)		2.81 (2.78)		
LSUN (church outdoor)		3.08 (2.94)		
CelebA		3.02 (2.97)		

Table 1: Bits/dim results for CIFAR-10, Imagenet, LSUN datasets and CelebA. Test results for CIFAR-10 and validation results for Imagenet, LSUN and CelebA (with training results in parenthesis for reference).

Real NVP: core idea

(real valued non-volume preserving transformations)

Real NVP: key property

(real valued non-volume preserving transformations)

$$\begin{cases} y_{1:d} = x_{1:d} \\ y & \text{The transformation is easily and} \\ analytically invertible! \end{pmatrix} + t(x_{1:d}) \\ \Leftrightarrow \begin{cases} x_{1:d} = y_{1:d} \\ x_{d+1:D} &= (y_{d+1:D} - t(y_{1:d})) \odot \exp\left(-s(y_{1:d})\right), \end{cases}$$

Real NVP: results

(real valued non-volume preserving transformations)

Dataset	PixelRNN [46]	Real NVP	Conv DRAW [22]	IAF-VAE [34]
CIFAR-10	3.00	3.49	< 3.59	< 3.28
Imagenet (32×32)	3.86 (3.83)	4.28 (4.26)	< 4.40 (4.35)	
Imagenet (64×64)	3.63 (3.57)	3.98 (3.75)	< 4.10 (4.04)	
LSUN (bedroom)		2.72 (2.70)		
LSUN (tower)		2.81 (2.78)		
LSUN (church outdoor)		3.08 (2.94)		
CelebA		3.02 (2.97)		

Real NVP: results

(real valued non-volume preserving transformations)

Density estimation using Real NVP https://arxiv.org/abs/1605.08803

Recurrent VAE: Iteratively Generating Data

DRAW: A Recurrent Neural Network For Image Generation https://arxiv.org/abs/1502.04623

Recurrent VAE: Iteratively Generating Data

Inference Model

Recurrent VAE: Iteratively Generating Imagenet

- Extend the process of generation and recognition to be sequential using recurrent neural networks.
- Inference and generation are convolutional maps.

Towards Conceptual Compression https://arxiv.org/abs/1604.08772

Recurrent VAE: Iteratively Generating Imagenet

Recurrent VAE

Pixel RNN

Applications of VAEs

Amortized inference with Normalizing Flows

MNIST

CIFAR10 patches

Making DNNs a little bit more Bayesian

Figure 3. Histogram of the trained weights of the neural network, for Dropout, plain SGD, and samples from Bayes by Backprop.

Weight Uncertainty in Neural Networks https://arxiv.org/abs/1505.05424

Making DNNs a little bit more Bayesian

Figure 2. Entropy gap ΔH_p (eq. (12)) between reversed and regular Penn Treebank test sets \times number of samples.

Bayesian Recurrent Neural Networks https://arxiv.org/abs/1704.02798

Improving semi-supervised with unsupervised features

Maintaining uncertainty allows us to:

- Speed up learning,
- Use less data, and
- Obtain better predictions.

Extracting relevant factors of variation

Early Visual Concept Learning with Unsupervised Deep Learning https://arxiv.org/abs/1606.05579

Discovering the dimensionality of the data

Stochastic Backpropagation and Approximate Inference in Deep Generative Models http://arxiv.org/pdf/1401.4082v3.pdf

Missing data imputation for 2D images

Stochastic Backpropagation and Approximate Inference in Deep Generative Models http://arxiv.org/pdf/1401.4082v3.pdf

Applications of Recurrent VAEs

Missing data imputation for 3D images

Unsupervised Learning of 3D Structure from Images https://arxiv.org/abs/1607.00662

8

E E E E E E E E E E E B B>RJBBBBBBBBBBBBB ZR 00 BALIN NKMEN NEODU Ee E 2 11 「ホホモリ用田田をきる」 Munn アホアモ ビリリシシ らてるし ppppp ドドゴミヨヨモをちろ DDDDD HUNNEEEE E bbbb エエエエ REPEREN pppp E E E E E N N N N ppppp

Lossy Compression

• Assuming we have access to the joint density

$$p(x,z) = p(x|z)\pi(z)$$

• A new sample **x** can be compressed by encoding it with the posterior density

- The number of bits necessary to communicate this density the KL-divergence $\frac{\mathrm{KL}(p(z|x);\pi(z))}{\ln(2)}$

Compression Rate

- Given the model
$$\ p(x,z) = p(x|z)\pi(z)$$

• Where
$$x \in \mathbb{R}^{d_x}$$

• We define compression rate as

$$r = \frac{1}{d} \mathbb{E}_{p(x)} \left[\frac{\mathrm{KL}(p(z|x); \pi(z))}{\ln(2)} \right]$$

Towards Conceptual Compression https://arxiv.org/abs/1604.08772

Understanding Images at multiple scales: How many bits are stored at each level?

$$\mathcal{F}(x) = \mathbf{E}_{\mathbf{q}}[-\ln \mathbf{p}(\mathbf{x}|\mathbf{z}) + \sum_{\mathbf{t}=1}^{\mathbf{L}} \mathrm{KL}(\mathbf{q}_{\mathbf{t}}(\mathbf{z}_{\mathbf{t}}|\mathbf{z}_{<\mathbf{t}},\mathbf{x});\mathbf{p}(\mathbf{z}_{\mathbf{t}}))]$$

$$r_{1} = \frac{1}{d} \mathbb{E}_{p(x)} \left[\mathrm{KL}(q_{1}(z_{1}|x));p(z_{1}) \right]$$

$$r_{2} = \frac{1}{d} \mathbb{E}_{p(x)} \left[\mathrm{KL}(q_{1}(z_{1}|x));p(z_{1}) \right] + \frac{1}{d} \mathbb{E}_{p(x)} \left[\mathrm{KL}(q_{2}(z_{2}|z_{<2},x));p(z_{2}) \right]$$

$$\vdots$$

$$r_{k} = \frac{1}{d} \sum_{k} \mathbb{E}_{p(x)} \left[\mathrm{KL}(q_{k}(z_{k}|z_{
Towards Conceptual Compression https://arxiv.org/abs/1604.08772$$

Understanding Images at multiple scales: How many bits are stored at each level?

Towards Conceptual Compression https://arxiv.org/abs/1604.08772

Original images

Compression rate: 0.05bits/dimension

Original images

Compression rate: 0.15bits/dimension

JPEG JPEG-2000 RVAE v1 RVAE v2

Original images

Compression rate: 0.2bits/dimension

JPEG JPEG-2000 RVAE v1

RVAE v2

Summary

• Variational methods in combination with deep learning are very powerful tools

for many problems in density estimation, information theory and control theory

- We have come a long way in scaling these ideas
- Several challenges remain:
 - Lower-variance gradient estimators
 - Several practical issues in training conditional models
 - Discrete variables
 - Applications to RL

A few references

- Stochastic Backpropagation and Approximate Inference in Deep Generative Models http://arxiv.org/abs/1401.4082
- Variational Inference with Normalizing Flows <u>https://arxiv.org/abs/1505.05770</u>
- Auto-Encoding Variational Bayes https://arxiv.org/abs/1312.6114
- Semi-Supervised Learning with Deep Generative Models http://arxiv.org/abs/1406.5298
- DRAW: A Recurrent Neural Network For Image Generation https://arxiv.org/abs/1502.04623
- Towards Conceptual Compression <u>https://arxiv.org/abs/1604.08772</u>
- Unsupervised Learning of 3D Structure from Images https://arxiv.org/abs/1607.00662
- One-Shot Generalization in Deep Generative Models http://arxiv.org/abs/1603.05106
- Normalizing Flows on Riemannian Manifolds https://arxiv.org/abs/1611.02304
- Improving Variational Inference with Inverse Autoregressive Flow https://arxiv.org/abs/1606.0493
- Weight Uncertainty in Neural Networks <u>https://arxiv.org/abs/1505.05424</u>
- Early Visual Concept Learning with Unsupervised Deep Learning https://arxiv.org/abs/1606.05579
- Bayesian Recurrent Neural Networks https://arxiv.org/abs/1704.02798
- Density estimation using Real NVP <u>https://arxiv.org/abs/1605.08803</u>

UAI 2017 Tutorial on Deep Generative Models

Shakir Mohamed and Danilo Rezende

Conference on Uncertainty in Artificial Intelligence Sydney, Australia August 11-15, 2017 Uai2017

Backup

Variational Methods for Channel Capacity Estimation

$$\mathcal{E}(\mathbf{s}) = \max_{\omega} \mathcal{I}^{\omega}(\mathbf{a}, \mathbf{s}' | \mathbf{s}) = \max_{\omega} \mathbb{E}_{p(s'|a, s)\omega(a|s)} \left[\log \left(\frac{p(\mathbf{a}, \mathbf{s}' | \mathbf{s})}{\omega(\mathbf{a}|\mathbf{s})p(\mathbf{s}'|\mathbf{s})} \right) \right]$$

Approximate Inference: Perturbative expansion compatible with variational bound

