

Auto-encoding Variational Bayes Kingma (2013)

how to perform efficient inference & learning in directed probabilistic model?

→ whilst in the presence of:

- continuous latent variables w/ intractable posterior distributions
 - large datasets

notes gen

A intractable posterior refers to a situation where the post prob dist cannot be computed in a simple or computationally viable way

Cannot:

- exact prob of specific para values
 - mean, var & sum stats
 - Samples from the dist likelihood prior

$$\text{recall Bayes: } P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D) \underset{\text{evidence}}{\text{Sum over poss. parameter values}}}$$

the main reason these problem become impossible
is because of the denominator (evidence $P(D)$)

- 1 no exact evidence to normalize the numerator and obtain posterior
(Den is constant & num is prop to Ths)
 - 2 high-Dims \rightarrow many paras / latent vars either imposs or comp exhaustion to calc / approx

(2)

Issues w/ Intractability:

① cannot sample from dist

② calc point estimates

③ compare between models (need $p(\theta)$)

Back to paper...

Paper introduces a stochastic variational inference & learning algo

↳ scales to large datasets

↳ works w/ intractable case

this approach is for models w/ cont latent variables and/or parameters

optimized w/ stochastic gradient des

i.i.d = independent & identically distributed

↳ RNNs & LSTMs handle non-i.i.d

~~background~~

this paper introduces the "Stochastic Gradient variational Bayes" (SGVB)

then a specific application called "Auto-encoding variational Bayes" for scenarios where dataset is i.i.d but w/ continuous latent variables + (common & import scenario)

Variational Bayes Optimizing of a Approx
to the Intractable Posterior (not poss)

Paper shows how re-para of variational
lower bounds yields a simple unbiased
estimator of the lower bound (SGVB)

↳ Allows eff approx posterior inference
in almost any model w/ cont latent
variables and f/pars

For 1.1.0 uses Auto-encode ~~VB~~ VB (AEVB)

- uses SGVB to make inference eff
- Approx post inference use simple ^{script} Aesmal to learn model pars
- avoid expensive methods such as Monte carlo

the learned approx posterior inf model
can be used for tasks: recognition,
denoise, representation & vis

~~Very~~ when NN is used for recognition
= Variational Auto encoder

SGVB is the foundational trick intro'd

Purpose is to solve prob of intractable exp
in variational inference

using re-para trick to get differentiable
unbiased estimator of variational lower bound

(4)

SGVB is applicable to all models w/ continuous latent variables

- regardless of whether data is i.i.d
- & regardless of how latent vars are structured across datapoints

AEVB

- Specialized SGVB for i.i.d dataset & continuous latent variables
- key innov of AEVB = intro of recognition model
- Recog model often implemented by NN
- learns to map datapoints (x) to the paras of - their approx posterior latent dist $q_\phi(z|x)$
- uses ancestral sample per datapoint to achieve — avoids using Monte Carlo
- if recognition NN = VAE architecture

Back to paper ...

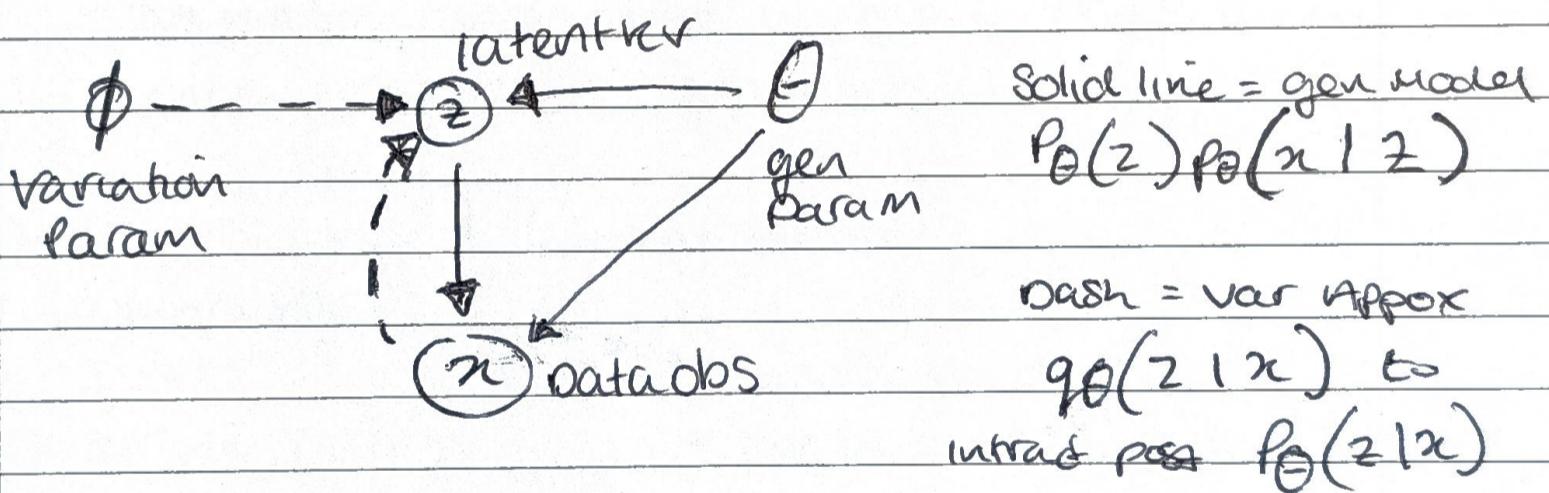
(5)

Method

Strategy presented can be used to derive a lower bound estimator (stochastic objective fun) for directed graphical model w/ cont latent vars

paper restricts section to common case w/ i.i.d dataset w/ latent vars

- ↳ Perform REML (ML) or max a posterior (MAP) inference on the global params
- ↳ & variational inference of latent vars



var & gen Params
learnt jointly

(6)

Problem Scenario

$$\text{Dataset } X = \{x^{(i)}\}_{i=1}^N \quad \text{i.i.d sample}$$

data is generated by random process involving an unobserved cont random variable z

Process of two steps:

$$z^{(i)}$$

(1) a value is generated from some prior distribution $P_{\theta}(z)$ (~~param~~ (needs params) θ defines shapes of dists)

(2) a value $x^{(i)}$ is generated from some conditional distribution $P_{\theta}(x|z)$

(latent)

z is an hidden, unobserved variable that fundamentally influences how x is generated

θ true param, $z^{(i)}$ latent variable

- (weights & bias)

- fixed

- underlying Distrs



- unobserved rvs

- diff values for diff data points

- often a specific latent var (or set) per data point

- inference to estimate (less) latents

Assume that prior $P_{\theta}(z)$ & likelihood $P_{\theta}(x|z)$ come from parametric distributions & that their PDFs are diff w.r.t both θ & z

both params θ & latents z are unknown to use

simplifies

important Assumps are not made about the marginal or posterior probs

Only interested in general algo that works efficient in a case of

- ① tractability
- ② large dataset - even batch is too costly
prefer small minibatch or single sample

Authors goal is to propose solutions to 3 probs:

1 ① Efficient Approx MC or MAP est for the params θ . these allow us to mimic the hidden random process & generate artificial data

2 ② EFF approx of Posterior inference of latent variable z given obs X & for choice of params θ

3 ③ EFF approx Marginal inference of the var X

to achieve all 3 the authors introduce the concept of a recognition model

$$q\phi(z|x)$$

this is an appox to the intractable true posterior $p\theta(z|x)$

(8)

later a method is introduced to learn ϕ & θ jointly

ϕ = recognition model params
 θ = generative model params

From a theory point, unobserved vars z are interpreted as latent repres or "Code"

the recognition Model $q_\phi(z|x)$ is a probabilistic encoder

given a datapoint x it produces a distribution over the possible values of the latent z

↳ from which x could have been generated

Also, $p_\theta(x|z)$ is a probabilistic Decoder as it takes z & produces x

2.2 Variational Bound

the marginal likelihood (Denominator)

↳ is composed of a sum over the marginal likelihoods of individual datapoints

$$\log p_{\theta}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}) = \sum_{i=1}^n \log p_{\theta}(\mathbf{x}^{(i)})$$

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

Divergence from true Post

Variational

lower Bound
on marg lik of
 $\mathbf{x}^{(i)}$

we want to find & optimize the lower bound
w.r.t variational params ϕ &
generative params θ

gradient of lower bkt ϕ is problematic

w/ Monte Carlo gradient exhibits very
high variance & is impractical for use

2.3 SFVB estimator & AEVB algo

Section introduces a practical estimator of the lower bound & its derivatives w.r.t the parameters

Assume posterior in the form $q_\phi(z|x)$

↳ technique can be applied to $q_\phi(z)$

* Bunch of stuff rearranging algos
↳ using Monte Carlo

2.4 Reparameterizing trick

Let z be a continuous random var

↳ $z \sim q_\phi(z|x)$ be a conditional distib

often able to express the random variable z as a deterministic variable $z = g_\phi(\epsilon, x)$

ϵ = Auxiliary variable w/ indep marginal $p(\epsilon)$

$g_\phi(\cdot)$ is some func parameterized by ϕ

repara is useful here as it can be used to rewrite an expectation w.r.t $q_\phi(z|x)$

such that the MC estimate is diff w.r.t ϕ

3. Variational Auto-encoder

Example uses as NN for the prob encoder
 $q_{\phi}(\epsilon | z)$

and where ϕ & θ are optimized jointly w/ the AEVB algo

more complicated

4. Related work

generative models are named as Decoders

Recognizing modes are labelled as encoders

hidden units number based on intricacy

3 latent variables

↳ says for high dim space results become unreliable