Variational Autoencoder Applications

Sargur N. Srihari srihari@cedar.buffalo.edu

Topics in VAE

- 1. VAE as a Generative Model
- https://cedar.buffalo.edu/~srihari/CSE676/21.1-VAE-Theory.pdf
- 1. VAE: The neural network perspective

https://cedar.buffalo.edu/~srihari/CSE676/21.2-VAE-NeuralNets.pdf

2. VAE Summary and Applications

VAE Summary and Applications

- 1. Summary of VAE
 - 1. Architecture
 - 2. Training
 - 3. Sample Generation
- 2. Deep Recurrent Attentive Writer (DRAW) Image Generation
- 3. Semi-Supervised Learning Semi-supervised learning
- 4. Interpolating between sentences
 Interpolating between sentences
- 5. Radiology

VAE Architecture



Source: https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html

Deep Learning

VAE Training (without/with reparameterization)



Trained by minimizing negative ELBO:

 $\boxed{l_i(\boldsymbol{\theta}, \boldsymbol{\phi}) = -E_{\boldsymbol{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z} \mid \boldsymbol{x}_i)} \left[\log p_{\boldsymbol{\theta}}(\boldsymbol{x}_i \mid \boldsymbol{z})\right] + KL\left(q_{\boldsymbol{\phi}}(\boldsymbol{z} \mid \boldsymbol{x}_i) \mid\mid p(\boldsymbol{z})\right)}$

5

Generating samples from a VAE

- Assume a trained VAE
- Generating a sample:
 - 1. Draw a sample from $p(\epsilon)$
 - 2. Run *z* through decoder to get $p(\boldsymbol{x}|\boldsymbol{z})$
 - It consists of parameters of a Bernoulli or multinomial
 - 3. Obtain sample from this distribution





Advantages/disadvantages of VAEs

- Advantages
 - Elegant
 - Generate from samples from N(0,1)
 - Theoretically pleasing
 - Maximizes ELBO in training
 - Excellent results
 - State-of-the-art for generative modeling
- Disadvantage
 - Samples from image VAEs tend to be blurry

Reasons for blurriness in VAE images

- **1. Effect of minimizing** $D_{KL}(p_{data}||p_{model})$
- Tend to assign high probability to points in data set but also to other points which are blurry

2. Gaussian $p_{\text{model}}(\boldsymbol{x}; g(\boldsymbol{z}))$

Maximizing lower bound is similar to training an autoencoder with mean squared error Has tendency to ignore features that occupy few pixels



Extending VAE to other architectures

- Extending VAE is straightforward
 - This a key advantage over Boltzmann machines
 - Which require careful design for tractability
- They work very well on a wide variety of differentiable operators
 - DRAW is a sophisticated recurrent encoder combined with an attention mechanism
 - Generation consists of visiting different small patches and drawing the values of the pixels at those points

Motivation for DRAW

- A person drawing a scene does it sequentially reassessing handiwork after each modification
 - Outlines are replaced by precise forms
 - lines are sharpened, darkened or erased
 - shapes are altered, and
 - the final picture emerges
- DRAW: Recurrent Neural Net For Image Generation
 - Google DeepMind 2015
- <u>https://www.youtube.com/watch?v=Zt-7MI9eKEo</u>



The Deep Recurrent Attentive Writer (DRAW) architecture

- Represents a natural form of image construction
- Parts of a scene are created independently from others,
- Approximate sketches are successively refined.



DRAW is a VAE

The encoder and decoder are both *RNNs* (LSTMs):

Approaches to image generation

- Most automatic image generation aim to generate entire scenes at once
 - This means that all pixels are conditioned on a single latent distribution
 - As well as precluding the possibility of iterative selfcorrection, the "one shot" approach is fundamentally difficult to scale to large images.
- DRAW attempts more natural image construction
 - Parts of a scene are created independently from others, and sketches are successively refined 12

A recurrent neural network (RNN)



13

Conventional VAE Design

Architecture



VAE is used in three different modes:

1. Generation:

a sample z is drawn from a prior p(z) and passed through the feedforward decoder network to compute the probability of the input p(x|z) given the sample

2. Inference:

the input x is passed to the encoder network, producing an approximate posterior q(z|x) over latent variables

3. Training:

z is sampled from q(z|x) and then used to compute the loss $KL[q(z|x)||p(z)]-\log p(x|z)$, which is minimized with SGD

The DRAW Network

- Basic network structure is similar to VAEs
 - Encoder network
 - Determines distribution over latent codes that capture salient information about input data
 - Decoder network
 - Receives samples from the code distribution and uses them to condition its own distribution over images
- There are three key differences
 - Described next

Differences between DRAW and VAE

- 1. Both encoder and de-coder are RNNs
 - A sequence of code samples is exchanged
 - Encoder is privy to decoder's previous outputs, allowing it to tailor the codes it sends according to the decoder's behavior so far
- 2. Decoder's outputs are successively added to distribution that ultimately generates the data as opposed to emitting distribution in a single step
- 3. Dynamically updated attention is used
 - to restrict both input region observed by encoder, and the output region modified by the decoder

Srihari

Encoder and decoder are both RNNs (LSTMs)

Deep Learning



- At each time-step, a sample z_t from prior $p(z_t)$ is passed to the recurrent decoder network, which then modifies part of the canvas matrix. The final canvas matrix c_T is used to compute $p(\boldsymbol{x}|\boldsymbol{z}_{1:T})$
- During inference the input is read at every time-step and the result is passed to the encoder RNN.

The RNNs at the previous time-step specify where to read. The output of the encoder RNN is used to compute the approximate posterior over the latent variables at that time-step.

Long Short Term Memory

- Explicitly designed to avoid the long-term dependency problem
- RNNs have the form of a repeating chain structure
 - <u>The repeating module has a simple structure such as tank</u>





- LSTMs also have a chain structure
 - but the repeating module has a different structure



Modeling the Drawing Mechanism

Iterative part

- Decoder's outputs are successively added to distribution that ultimately generates the data, as opposed to emitting it in a single step
- Attention mechanism
 - To model working on one part of the image, and then another
 - Used to restrict both input region observed by encoder, and output region modified by the decoder
- Network decides at each time-step:
 - where to read, where to write and what to write

Attention Mechanism



A classic image captioning system

would encode the image, using a pre-trained CNN that would produce a hidden state

With an attention mechanism,

the image is first divided into parts, and we compute with a CNN representations of each part $h_1, \dots h_n$

When the RNN is generating a new word, the attention mechanism is focusing on the relevant part of the image,

so the decoder only uses specific parts of the image

Selective Attention Model

- N×N grid of Gaussian filters
 - Grid center (g_X, g_Y) and stride δ determine mean of filter at patch (i,j): $\mu^i_X = g_X + (i - N/2 - 0.5)\delta$, $\mu^j_Y = g_Y + (j - N/2 - 0.5)\delta$

3x 3 grid of filters





- Three $N \times N$ patches extracted from the image (N=12). Green rectangles indicate the boundary and precision (σ) of the patches, while the patches themselves are shown to the right.
 - 1. Top patch has a small δ and high σ , giving a zoomed-in but blurry view of the center of the digit;
 - 2. Middle patch has large δ and low σ , effectively down sampling the whole image;
 - 3. Bottom patch has high δ and σ

VAE and Semi-supervised Learning

- Semi-supervised learning falls in between unsupervised and supervised learning
 - make use of both labeled and unlabeled data in supervised learning (e.g. classification, regression)
- Billions of unlabeled images on the internet, only a small fraction is labeled

– Humans can identify anteaters after seeing very few

Goal is to get the best performance from a tiny data set

Deep Learning

VAE for semi-supervised learning



PGM for semi-supervised VAE

Generative model P

Inference model Q



The incoming joint connections to each variable are deep neural networks with parameters θ and ϕ

a are auxiliary variables such that: q(a,z|x) = q(z|a,x)q(a|x) and the marginal distribution q(z|x) can fit more complicated posteriors p(z|x)

Word-to-Vec for Language

Word-to-Vec:

e.g., sentence Have a great day Start with one hot-vector of size |V| $V = \{$ have, a, great, day $\}$

Common Bag of Words

The word great is input to predict day



Result is an embedding of great in vector space of N dimensions

With C context words

Have a great is input to predict day



Skip-gram model Predicting the context

For each context position, we get *C* probability distributions of *V* probabilities



Sentence prediction by conventional autoencoder

Sentences produced by greedily decoding from points between two sentence encodings with a conventional autoencoder The intermediate sentences are not plausible English

i went to the store to buy some groceries . *i store to buy some groceries* . *i were to buy any groceries* . *horses are to buy any groceries* . *horses are to buy any animal* . *horses the favorite any animal* . *horses the favorite favorite animal* . **horses are my favorite animal** .

VAE Language Model

 Words are represented using a learned dictionary of embedding words



VAE sentence interpolation

- Paths between random points in VAE space
- Intermediate sentences are grammatical
- Topic and syntactic structure are consistent

"i want to talk to you ." "i want to be with you ." "i do n't want to be with you ." i do n't want to be with you . she did n't want to be with him .

he was silent for a long moment . he was silent for a moment . it was quiet for a moment . it was dark and cold . there was a pause . it was my turn .

VAE for Radiology

Combines two types of models:

- discriminative & generative models



Models trained jointly using variational EM framework

Left: Discriminative deep nn model Input: observed variables Generates posterior distributions over latent variables and possibly (if unobserved) class labels.

Performs Inference of latent variables necessary to perform variational updates

Right: generative PGM with inputs:

- 1. class label y (diseases)
- 2. nuisance variables s (hospital identifiers)
- 3. latent variables z (size, shape, other

brain properties)

Provides causality of observation