

Lecture 8: Deep Generative Models Efstratios Gavves

UVA DEEP LEARNING COURSE – EFSTRATIOS GAVVES

- o Early Generative Models
- Restricted Boltzmann Machines
- o Deep Boltzmann Machines
- Deep Belief Network
- o Contrastive Divergence
- o Gentle intro to Bayesian Modelling and Variational Inference
- Variational Autoencoders
- Normalizing Flows

Plug in the model density function to likelihood
Then maximize the likelihood

Problems

- Design complex enough model that meets data complexity
- At the same time, make sure model is computationally tractable
- More details in the next lecture



Restricted Boltzmann Machines Deep Boltzmann Machines Deep Belief Nets



• We can define an explicit density function over all possible relations ψ_c between the input variables x_c

$$p(x) = \prod_{c} \psi_{c} \left(x_{c} \right)$$

• Quite inefficient \rightarrow think of all possible relations between $256 \times 256 = 65K$ input variables

• Not just pairwise

• Solution: Define an energy function to model these relations

• First, define an energy function -E(x) that models the joint distribution

$$p(x) = \frac{1}{Z} \exp(-E(x))$$

 $\circ Z$ is a normalizing constant that makes sure p(x) is a pdf: $\int p(x) = 1$

$$Z = \sum_{x} \exp(-E(x))$$

• Well understood in physics, mathematics and mechanics

• A Boltzmann distribution (also called Gibbs distribution) is a probability distribution, probability measure, or frequency distribution of particles in a system over various possible states

• The distribution is expressed in the form

$$F(state) \propto \exp(-\frac{E}{kT})$$

 $\circ E$ is the state energy, k is the Boltzmann constant, T is the thermodynamic temperature

https://en.wikipedia.org/wiki/Boltzmann_distribution

Problem with Boltzmann Distribution?

 \circ Initially, thought of in the context of binary variables x

 Assuming binary variables x the normalizing constant has very high computational complexity

 \circ For *n*-dimensional *x* we must enumerate all possible 2^n operations for *Z*

• Clearly, gets out of hand for any decent *n*

• Solution: Consider only pairwise relations

• The energy function becomes

$$E(x) = -x^T W x - b^T x$$

 $\circ x$ is considered binary

 $\circ x^T W x$ captures correlations between input variables

$o b^T x$ captures the model prior

• The energy that each of the input variable contributes itself

Problem with Boltzmann Machines?

• Still too complex and high-dimensional

- o If x has $256 \times 256 = 65536$ dimensions
- \circ The pairwise relations need a huge W: 4.2 billion dimensions

<u>Just for connecting two layers!</u>

• Solution: Consider latent variables for model correlations

• Restrict the model energy function further to a bottleneck over latents h

$$E(x) = -x^T W h - b^T x - c^T h$$

$$\mathbf{O}E(x) = -x^T W h - b^T x - c^T h$$

• The $x^T W h$ models correlations between x and the latent activations via the parameter matrix W

• The $b^T x$, $c^T h$ model the priors

• Restricted Boltzmann Machines (RBM) assume x, h to be binary

• Energy function:
$$E(x) = -x^T W h - b^T x - c^T h$$

 $p(x) = \frac{1}{Z} \sum_{h} \exp(-E(x,h))$
• Not in the form $\propto \exp(x)/Z$ because of the Σ



• The F(x) defines a bipartite graph with undirected connections • Information flows forward and backward



• The hidden units h_j are independent to each other conditioned on the visible units

$$p(h|x) = \prod_{j} p(h_{j}|x,\theta)$$

• The visible units x_i are independent to each other conditioned on the hidden units

$$p(x|h) = \prod_{i} p(x_i|h,\theta)$$



Training RBMs

• The conditional probabilities are defined as sigmoids $p(h_j | x, \theta) = \sigma(W_{\cdot j} x + b_j)$ $p(x_i | h, \theta) = \sigma(W_{\cdot i} x + c_i)$

o Maximize log-likelihood

$$p(x) = \frac{1}{Z} \exp(-F(x))$$

 $\mathcal{L}(\theta) = \frac{1}{N} \sum \log p(x_n | \theta)$



DEEP GENERATIVE MODELS - 18







• Let's take the gradients

$$\frac{\partial \log p(x_n|\theta)}{\partial \theta} = -\frac{\partial F(x_n)}{\partial \theta} - \frac{\partial \log Z}{\partial \theta} \\ = -\sum_{h}^{\partial \theta} p(h|x_n, \theta) \frac{\partial E(x_n|h, \theta)}{\partial \theta} + \sum_{\tilde{x}, h}^{\partial \theta} p(\tilde{x}, h|\theta) \frac{\partial E(\tilde{x}, h|\theta)}{\partial \theta}$$

Hidden unit (features)



Easy because we just substitute in the definitions the x_n and sum over h
 Hard because you need to sum over both x̃, h which can be huge
 It requires approximate inference, e.g., MCMC

• Approximate the gradient with Contrastive Divergence

• Specifically, apply Gibbs sampler for k steps and approximate the gradient $\partial \log p(x_n|\theta) \qquad \partial E(x_n, h_0|\theta) \quad \partial E(x_k, h_k|\theta)$ $\partial \theta$ $\partial \theta$ $\partial \theta$ $\mathbf{h}_0 \sim \mathbf{P}(\mathbf{h}|\mathbf{x})$ $\mathbf{h}_1 \sim \mathbf{P}(\mathbf{h}|\mathbf{x}_1)$ Observations Reconstructions $\mathbf{x}_1 \sim \mathbf{P}(\mathbf{x}|\mathbf{h})$ Х

Hinton, Training Products of Experts by Minimizing Contrastive Divergence, Neural Computation, 2002

o RBMs are just one layer

• Use RBM as a building block

• Stack multiple RBMs one on top of the other $p(x, h_1, h_2) = p(x|h_1) \cdot p(h_1|h_2)$

- Deep Belief Networks (DBN) are directed models
 - The layers are densely connected and have a single forward flow
 - This is because the RBM is directional, $p(x_i|h, \theta) = \sigma(W_i x + c_i)$, namely the input argument has only variable only from below



- Stacking layers again, but now with connection from the above and from the below layers
- Since it's a Boltzmann machine, we need an energy function

$$E(x, h_1, h_2 | \theta) = x^T W_1 h_1 + h_1^T W_2 h_2 + h_2^T W_3 h_3$$
$$p(h_2^k | h_1, h_3) = \sigma(\sum_j W_1^{jk} h_1^j + \sum_l W_3^{kl} h_3^k)$$



• Schematically similar to Deep Belief Networks

- But, Deep Boltzmann Machines (DBM) are undirected models
 - Belong to the Markov Random Field family
- So, two types of relationships: bottom-up and upbottom

$$p(h_2^k | h_1, h_3) = \sigma(\sum_j W_1^{jk} h_1^j + \sum_l W_3^{kl} h_3^k) \quad \mathbf{I}$$



• Computing gradients is intractable

o Instead, variational methods (mean-field) or sampling methods are used



Variational Inference



Some (probabilistic) terminology

 \circ Observed variables x

- $_{
 m O}$ Latent variables heta
- Both unobservable model parameters w and unobservable model activations z

 $\circ \theta = \{w, z\}$

• Joint probability density function (pdf): $p(x, \theta)$

• Marginal pdf: $p(x) = \int_{\theta} p(x, \theta) d\theta$

- Prior pdf \rightarrow marginal over input: $p(\theta) = \int_x p(x, \theta) dx$
- Usually a user defined pdf
- Posterior pdf: $p(\theta|x)$
- o Likelihood pdf: $p(x|\theta)$







• Conjugate priors

 when posterior and prior belong to the same family, so the joint pdf is easy to compute

- Point estimate approximations of latent variables
 - instead of computing a distribution over all possible values for the variable
 - compute one point only
 - e.g. the most likely (maximum likelihood or max a posteriori estimate)

 $\theta^* = \arg_{\theta} \max p(x|\theta)p(\theta) (MAP)$ $\theta^* = \arg_{\theta} \max p(x|\theta) (MLE)$

 Quite good when the posterior distribution is peaky (low variance)



• Estimate the posterior density p(θ|x) for your training data x
• To do so, need to define the prior p(θ) and likelihood p(x|θ) distributions
• Once the p(θ|x) density is estimated, Bayesian Inference is possible
• p(θ|x) is a (density) function, not just a single number (point estimate)

• But how to estimate the posterior density?

- Markov Chain Monte Carlo (MCMC) \rightarrow Simulation-like estimation
- $^{\circ}$ Variational Inference \rightarrow Turn estimation to optimization

• Estimating the true posterior $p(\theta|x)$ is not always possible • especially for complicated models like neural networks

- Variational Inference assumes another function $q(\theta|\varphi)$ with which we want to approximate the true posterior $p(\theta|x)$
 - $\circ q(\theta|\varphi)$ is the approximate posterior
 - $^{\rm o}$ Note that the approximate posterior does not depend on the observable variables x

• We approximate by minimizing the **reverse** KL-divergence w.r.t. $\varphi^* = \arg\min_{\varphi} KL(q(\theta|\varphi)||p(\theta|x))$

• Turn inference into optimization









Variational Inference - Evidence Lower Bound (ELBO)

o Given latent variables θ and the approximate posterior

$$q_{\varphi}(\theta) = q(\theta|\varphi)$$

• What about the log marginal $\log p(x)$?

Variational Inference - Evidence Lower Bound (ELBO)

 \circ Given latent variables heta and the approximate posterior

$$q_{\varphi}(\theta) = q(\theta|\varphi)$$

• We want to maximize the marginal p(x) (or the log marginal $\log p(x)$

$$\log p(x) \ge \mathbb{E}_{q_{\varphi}(\theta)} \left[\log \frac{p(x,\theta)}{q_{\varphi}(\theta)} \right]$$

Evidence Lower Bound (ELBO): Derivations



• The log marginal is



$$\begin{split} KL\left[q(Z)\||p(Z|X)\right] &= \int_{Z} q(Z)\log\frac{q(Z)}{p(Z|X)} \\ &= -\int_{Z} q(Z)\log\frac{p(Z|X)}{q(Z)} \\ &= -\left(\int_{Z} q(Z)\log\frac{p(X,Z)}{q(Z)} - \int_{Z} q(Z)\log p(X)\right) \\ &= -\int_{Z} q(Z)\log\frac{p(X,Z)}{q(Z)} + \log p(X)\int_{Z} q(Z) \\ &= -L + \log p(X) \end{split}$$

UVA DEEP LEARNING COURSE – EFSTRATIOS GAVVES

$$\geq \mathbb{E}_{q_{\varphi}(\theta)} \left[\log \frac{p(x,\theta)}{q_{\varphi}(\theta)} \right]$$

$$= \mathbb{E}_{q_{\varphi}(\theta)} [\log p(x|\theta)] + \mathbb{E}_{q_{\varphi}(\theta)} [\log p(\theta)] - \mathbb{E}_{q_{\varphi}(\theta)} [\log q_{\varphi}(\theta)]$$

$$= \mathbb{E}_{q_{\varphi}(\theta)} [\log p(x|\theta)] - \mathrm{KL}(q_{\varphi}(\theta)||p(\theta))$$

$$= \mathrm{ELBO}_{\theta,\varphi}(x)$$

o Maximize reconstruction accuracy $\mathbb{E}_{q_{\varphi}(\theta)}[\log p(x|\theta)]$

• While minimizing the KL distance between the prior $p(\theta)$ and the approximate posterior $q_{\varphi}(\theta)$

ELBO: Formulation 2

$$\geq \mathbb{E}_{q_{\varphi}(\theta)} \left[\log \frac{p(x,\theta)}{q_{\varphi}(\theta)} \right]$$

= $\mathbb{E}_{q_{\varphi}(\theta)} [\log p(x,\theta)] - \mathbb{E}_{q_{\varphi}(\theta)} [\log q_{\varphi}(\theta)]$
= $\mathbb{E}_{q_{\varphi}(\theta)} [\log p(x,\theta)] + H(\theta)$
= $\mathrm{ELBO}_{\theta,\varphi}(x)$

• Maximize something like negative Boltzmann energy $\mathbb{E}_{q_{\varphi}(\theta)}[\log p(x, \theta)]$ • While maximizing the entropy the approximate posterior $q_{\varphi}(\theta)$ • Avoid collapsing latents θ to a single value (like for MAP estimates) o It is easy to see that the ELBO is directly related to the marginal

$$\log p(x) = \text{ELBO}_{\theta,\varphi}(x) + KL(q_{\varphi}(\theta)||p(\theta|x))$$

 \circ You can also see $ELBO_{\theta,\phi}(x)$ as Variational Free Energy

o It is easy to see that the ELBO is directly related to the marginal $ELBO_{\theta,\phi}(x) =$

o It is easy to see that the ELBO is directly related to the marginal $ELBO_{\theta,\omega}(x) =$ $= \mathbb{E}_{q_{\varphi}(\theta)}[\log p(x,\theta)] - \mathbb{E}_{q_{\varphi}(\theta)}[\log q_{\varphi}(\theta)]$ $= \mathbb{E}_{q_{\varphi}(\theta)}[\log p(\theta|x)] + \mathbb{E}_{q_{\varphi}(\theta)}[\log p(x)] - \mathbb{E}_{q_{\varphi}(\theta)}[\log q_{\varphi}(\theta)]$ $= \mathbb{E}_{q_{\varphi}(\theta)}[\log p(x)] - KL(q_{\varphi}(\theta)||p(\theta|x))$ $= \log p(x) - KL(q_{\varphi}(\theta) || p(\theta | x))$ $\log p(x)$ does not depend on $q_{arphi}(heta)$ $\mathbb{E}_{q_{\varphi}(\theta)}[1]=1$ \Rightarrow $\log p(x) = \text{ELBO}_{\theta, \varphi}(x) + KL(q_{\varphi}(\theta)||p(\theta|x))$ • You can also see $ELBO_{\theta, \omega}(x)$ as Variational Free Energy

 $o \log p(x) = \text{ELBO}_{\theta,\varphi}(x) + KL(q_{\varphi}(\theta)||p(\theta|x))$

• The log-likelihood $\log p(x)$ constant \rightarrow does not depend on any parameter • Also, $\text{ELBO}_{\theta,\varphi}(x) > 0$ and $KL(q_{\varphi}(\theta)||p(\theta|x)) > 0$

- 1. The higher the Variational Lower Bound $\text{ELBO}_{\theta,\varphi}(\mathbf{x})$, the smaller the difference between the approximate posterior $q_{\varphi}(\theta)$ and the true posterior $p(\theta|\mathbf{x}) \rightarrow$ better latent representation
- 2. The Variational Lower Bound $ELBO_{\theta,\varphi}(x)$ approaches the log-likelihood \rightarrow better density model

- The variational distribution $q(\theta|\varphi)$ does not depend directly on data • Only indirectly, via minimizing its distance to the true posterior $KL(q(\theta|\varphi)||p(\theta|x))$
- \circ So, with $q(\theta|\varphi)$ we have a major optimization problem
- The approximate posterior must approximate the whole dataset $x = [x_1, x_2, ..., x_N]$ jointly
- \circ Different neural network weights for each data point x_i

 Better share weights and "amortize" optimization between individual data points

$$q(\theta|\varphi) = q_{\varphi}(\theta|x)$$

 \circ Predict model parameters heta using a arphi-parameterized model of the input x

- Use amortization for data-dependent parameters that depend on data
- E.g., the latent activations that are the output of a neural network layer: $z \sim q_{\varphi}(z|x)$

- The original view on Variational Inference is that $q(\theta|\varphi)$ describes the approximate posterior of the dataset as a whole
- Imagine you don't want to make a practical model that returns latent activations for a specific input
- Instead, you want to optimally approximate the true posterior of the unknown weights with an model with latent parameters
- $_{\rm O}$ It doesn't matter if these parameters are "latent activations" z or "model variables" w

• Let's rewrite the ELBO a bit more explicitly $ELBO_{\theta,\varphi}(x) = \mathbb{E}_{q_{\varphi}(\theta)}[\log p(x|\theta)] - KL(q_{\varphi}(\theta)||p(\theta))$ $= \mathbb{E}_{q_{\varphi}(z|x)}[\log p_{\theta}(x|z)] - KL(q_{\varphi}(z|x)||p_{\lambda}(z))$

$\circ p_{\theta}(x|z)$ instead of $p(x|\theta)$

o I.e., the likelihood model $p_{\theta}(x|z)$ has weights parameterized by θ o Conditioned on latent model activations parameterized by z

• Let's rewrite the ELBO a bit more explicitly $ELBO_{\theta,\varphi}(x) = \mathbb{E}_{q_{\varphi}(\theta)}[\log p(x|\theta)] - KL(q_{\varphi}(\theta)||p(\theta))$ $= \mathbb{E}_{q_{\varphi}(z|x)}[\log p_{\theta}(x|z)] - KL(q_{\varphi}(z|x)||p_{\lambda}(z))$

$\circ p_{\lambda}(z)$ instead of $p(\theta)$

o I.e., a λ -parameterized prior only on the latent activations zo <u>Not on model weights</u>

• Let's rewrite the ELBO a bit more explicitly $ELBO_{\theta,\varphi}(x) = \mathbb{E}_{q_{\varphi}(\theta)}[\log p(x|\theta)] - KL(q_{\varphi}(\theta)||p(\theta))$ $= \mathbb{E}_{q_{\varphi}(z|x)}[\log p_{\theta}(x|z)] - KL(q_{\varphi}(z|x)||p_{\lambda}(z))$

 $\circ q_{\varphi}(z|x)$ instead of $q(\theta|\varphi)$

o The model $q_{\varphi}(z|x)$ approximates the posterior density of the latents zo The model weights are parameterized by φ • ELBO_{θ,φ} $(x) = \mathbb{E}_{q_{\varphi}(z|x)}[\log p_{\theta}(x|z)] - \text{KL}(q_{\varphi}(z|x)||p_{\lambda}(z))$ • How to model $p_{\theta}(x|z)$ and $q_{\varphi}(z|x)$? $\circ \text{ELBO}_{\theta,\varphi}(x) = \mathbb{E}_{q_{\varphi}(z|x)}[\log p_{\theta}(x|z)] - \text{KL}(q_{\varphi}(z|x)||p_{\lambda}(z))$ $\circ \text{How to model } p_{\theta}(x|z) \text{ and } q_{\varphi}(z|x)?$

• What about modelling them as neural networks

Variational Autoencoders

 \circ Input x is an image

infers/recognizes the latent codes

get the Vanilla Autoencoder

• The approximate posterior $q_{\varphi}(z|x)$ is a CovnNet (or MLP)

• Also known as encoder or inference or recognition network, because it

• The likelihood density $p_{\theta}(x|z)$ is an inverted ConvNet (or MLP) $p_{\lambda}(z)$

o If we ignore the distribution of the latents $z,\,p_\lambda(z))$, then we network

 \circ Given input the output is a feature map from a latent variable z

 \circ Given the latent z as input, it reconstructs the input \widetilde{x}

• Also known as decoder or generator network

 $p_{\theta}(x|z)$

 $q_{\varphi}(z|x)$

Decoder/Generator

network

• Maximize the Evidence Lower Bound (ELBO)

• Or minimize the negative ELBO

 $\mathcal{L}(\theta, \varphi) = \mathbb{E}_{q_{\varphi}(z|x)}[\log p_{\theta}(x|z)] - \mathrm{KL}(q_{\varphi}(z|x)||p_{\lambda}(z))$

 \odot How to we optimize the ELBO?

• Maximize the Evidence Lower Bound (ELBO)

• Or minimize the negative ELBO

$$\mathcal{L}(\theta,\varphi) = \mathbb{E}_{q_{\varphi}(Z|x)}[\log p_{\theta}(x|Z)] - \mathrm{KL}(q_{\varphi}(Z|x)||p_{\lambda}(Z))$$
$$= \int_{Z} q_{\varphi}(z|x) \log p_{\theta}(x|z) \, dz - \int_{Z} q_{\varphi}(z|x) \log \frac{q_{\varphi}(z|x)}{p_{\lambda}(z)} \, dz$$

 $_{\odot}$ Forward propagation ightarrow compute the two terms

- The first term is an integral (expectation) that we cannot solve analytically. So, we need to sample from the pdf instead
 - When $p_{\theta}(x|z)$ contains even a few nonlinearities, like in a neural network, the integral is hard to compute analytically

• Maximize the Evidence Lower Bound (ELBO)

• Or minimize the negative ELBO

$$\mathcal{L}(\theta,\varphi) = \mathbb{E}_{q_{\varphi}(Z|x)}[\log p_{\theta}(x|Z)] - \mathrm{KL}(q_{\varphi}(Z|x)||p_{\lambda}(Z))$$
$$= \int_{Z} q_{\varphi}(z|x) \log p_{\theta}(x|z) \, dz - \int_{Z} q_{\varphi}(z|x) \log \frac{q_{\varphi}(z|x)}{p_{\lambda}(z)} \, dz$$

 \circ Forward propagation \rightarrow compute the two terms

- The first term is an integral (expectation) that we cannot solve analytically.
- When $p_{\theta}(x|z)$ contains even a few nonlinearities, like in a neural network, the integral is hard to compute analytically
- So, we need to sample from the pdf instead
- VAE is a stochastic model

• The second term is the KL divergence between two distributions that we know

$\circ \int_{z} q_{\varphi}(z|x) \log p_{\theta}(x|z) \, dz$

- The first term is an integral (expectation) that we cannot solve analytically.
- When $p_{\theta}(x|z)$ contains even a few nonlinearities, like in a neural network, the integral is hard to compute analytically
- As we cannot compute analytically, we sample from the pdf instead
 - Using the density $q_{\varphi}(z|x)$ to draw samples
 - $^{\circ}$ Usually one sample is enough ightarrow stochasticity reduces overfitting
- VAE is a stochastic model
- The second term is the KL divergence between two distributions that we know

$$\circ \int_{z} q_{\varphi}(z|x) \log \frac{q_{\varphi}(z|x)}{p_{\lambda}(z)} dz$$

- The second term is the KL divergence between two distributions that we know
- \circ E.g., compute the KL divergence between a centered N(0,1) and a noncentered $N(\mu,\sigma)$ gaussian

• We set the prior $p_{\lambda}(z)$ to be the unit Gaussian $p(z) \sim N(0, 1)$

We set the likelihood to be a Bernoulli for binary data

$p(x|z) \sim Bernoulli(\pi)$

 $_{\rm O}$ We set $q_{\varphi}({\bf z}|{\bf x})$ to be a neural network (MLP, ConvNet), which maps an input ${\bf x}$ to the Gaussian distribution, specifically it's mean and variance

$$^{\circ}\mu_{z}$$
, $\sigma_{z} \sim q_{\varphi}(\mathbf{z}|\mathbf{x})$

 $^{\circ}$ The neural network has two outputs, one is the mean μ_{χ} and the other the $\sigma_{\chi},$ which corresponds to the covariance of the Gaussian



• We set $p_{\theta}(\mathbf{x}|\mathbf{z})$ to be an inverse neural network, which maps Z to the Bernoulli distribution if our outputs binary (e.g. Binary MNIST)

• Good exercise: Derive the ELBO for the standard VAE



VAE: Interpolation in the latent space



Summary

UVA DEEP LEARNING COURSE EFSTRATIOS GAVVES DEEP GENERATIVE MODELS - 71 • Gentle intro to Bayesian Modelling and Variational Inference • Restricted Boltzmann Machines • Deep Boltzmann Machines O Deep Belief Network Contrastive Divergence • Variational Autoencoders • Normalizing Flows