

# Deepmind - Modern LVM

Saturday, 27 February 2021 11:48

- Generative models are probabilistic models of high-dimensional data.

- Describe the probabilistic process of generating an observation.
  - Emphasize on capturing dependence between the dimensions.
  - Provide a way of generating new data points.

- Types of generative models:

## 1. Autoregressive models:

- RNN & transformer language models, MADE, Pixel CNN, WaveNet

## 2. Latent variable models:

- Tractable: Invertible / flow-based models (RealNVP, Glow...)

- Intractable: Variational autoencoders

## 3. Implicit models:

- Generative Adversarial Networks (GANs)

- Auto-regressive models:

- Model the 1-dimensional conditional distributions instead of modeling the joint distribution directly;

$$\text{Based on the chain rule } p(x) = \prod_{d=1}^D p(x_d | x_1, \dots, x_{d-1})$$

- Trained with maximum likelihood.

- Pros: 1. 1-dimensional dist. are easy to model.

- 2. Simple and efficient training.

- 3. No sampling at training time.

- Cons: 1. Slow, sequential generation (one dimension at a time)

- 2. Usually much better at modeling local structure than global structure.

- Latent variable models:

- Specify the generative process in terms of unobserved/latent variables and the transformation that maps them to the observation.

- Trained with maximum likelihood (usually with approximations)

- Pros: 1. Powerful and well-understood framework.

- 2. Easy to incorporate prior knowledge / structure into models well-suited for representation learning / interpretability

- 3. Fast generation.

- Cons: 1. Conceptually more complex than fully observed models.

- 2. "Mode collapse": Models ignore regions of the data dist.

- 3. Training can be unstable and require many tricks to work.

- Generative Adversarial Networks:

- Model: NN that maps noise vectors to observations.

- Training: use the learning signal from a classifier trained to discriminate between samples from the model and the training data.

- Pros: 1. Can generate very realistic images.

- 2. Conceptually simple implementation.

- 3. Fast generation.

- Cons: 1. Cannot be used to compute probability of observations.

- 2. "Mode collapse": Models ignore regions of the data dist.

- 3. Training can be unstable and require many tricks to work.

- Latent variable models: define a distribution over observations  $x$  by using a (vector) latent variable  $z$  and specifying:

- The prior distribution  $p(z)$  for the latent variable.

- The likelihood  $p(x|z)$  connects the latent variable to the observation.

- The prior and the likelihood define the joint distribution  $p(x,z) = p(x|z)p(z)$

- We need the marginal likelihood  $p(x)$  and the posterior  $p(z|x)$

- We can think of latent variable value as an explanation for the observation.

- To generate an observation from the model:  $z \sim p(z) \Rightarrow x \sim p(x|z)$

- Inference: Going from observation to latent values.

- i.e. computing the posterior distribution for the given observation:

$$p(z|x) = \frac{p(x,z)}{\int p(x,z) dz}$$

This requires computing the marginal likelihood:  $p(x) = \int p(x,z) dz$

- Generating pairs  $(x,z)$ :  $z \sim p(z) \Rightarrow x \sim p(x|z)$

$$\text{or } x \sim p(x) \Rightarrow z \sim p(z|x)$$

- Maximizing likelihood:  $\theta^* = \arg \max \sum_{i=1}^n \log p_\theta(x_i)$

- The gradient for a single point:

$$\nabla_\theta \log p_\theta(x) = \int p_\theta(z|x) \nabla_\theta \log p_\theta(z) dz$$

- Invertible models (aka. normalizing flows):

- Approximate the data distribution by transforming the prior distribution using an invertible function; can be done  $\otimes$  jointly for all dimensions or autoregressively; one dimension at a time.

- Generate observations by applying an invertible and differentiable transformation  $f_\theta(z)$  to sample from the prior:  $z \sim p(z) \Rightarrow x = f_\theta(z)$

The marginal likelihood: we can relate densities of  $x$  and  $z$ :

$$p_\theta(x) = p_\theta(z) \left| \det \frac{\partial z}{\partial x} \right| = p_\theta(f_\theta^{-1}(x)) \left| \det \frac{\partial f_\theta^{-1}(x)}{\partial x} \right|$$

The determinant of the Jacobian takes into account the local change in volume when applying the transformation.

- Independent Component Analysis (ICA):

- uses fractional prior and a linear mapping (a square matrix):

$$z \sim \prod_i p(z_i) \Rightarrow x = Az$$

- Limitations: 1. The latent space and data space dimensionality must be the same.

- 2. The latent variable have to be continuous.

- 3. Observations have to be continuous or quantized.

- 4. Expressive models require a lot of layers, parameters.

- 5. Lack of flexibility in model design.

- Solution: use approximate inference for intractable models.

- Approximate inference:

- 1. Markov Chain Monte Carlo: generate samples from the exact posterior using a Markov Chain:

- +: very general; exact in the limit of infinite time/computation.

- : computationally expensive; convergence is hard to diagnose.

- 2. Variational inference: approximate the posterior with a tractable distribution; e.g. fully factorized or auto-regressive.

- +: fairly efficient, inference is reduced to optimization w.r.t. dist. parameters.

- : cannot trade computation for accuracy easily.

- Variational inference:

- We approximate the exact posterior  $p_\theta(z|x)$  with a variational posterior  $q_\phi(z|x)$

- We can use any distribution for  $q_\phi(z|x)$  as long as:

- 1. We can sample from it.

- 2. We can compute  $\log q_\phi(z|x)$  and its gradient w.r.t.  $\phi$

Training: The variational posterior induces a variational lower bound  $L_{\theta,\phi}(x)$  on the marginal log-likelihood.

→ we train the model by maximizing  $L_{\theta,\phi}(x)$  w.r.t both

model parameters  $\theta$  and the variational parameter  $\phi$ .

Bounding the marginal log-likelihood:

$$\log p_\theta(x) = \log \int p_\theta(x,z) dz = \log \int q_\phi(z) \frac{p_\theta(x,z)}{q_\phi(z)} dz \geq \int q_\phi(z) \log \frac{p_\theta(x,z)}{q_\phi(z)} dz$$

$$\Rightarrow L_{\theta,\phi}(x) = \mathbb{E}_{q_\phi(z)} [\log \frac{p_\theta(x,z)}{q_\phi(z)}] \leq \log p_\theta(x)$$

- Maximizing likelihood:  $\theta^* = \arg \max \sum_{i=1}^n \log p_\theta(x_i)$

- The gradient for a single point:

$$\nabla_\theta \log p_\theta(x) = \int p_\theta(z|x) \nabla_\theta \log p_\theta(z) dz$$

- Evidence Lower Bound (ELBO):

- Is obtained by using the variational posterior  $q_\phi(z|x)$  as  $q_\phi(z)$

$$L_{\theta,\phi}(x) = \mathbb{E}_{q_\phi(z|x)} [\log \frac{p_\theta(x,z)}{q_\phi(z)}]$$

- Limitations: 1. The latent space and data space dimensionality must be the same.

- 2. The latent variable have to be continuous.

- 3. Observations have to be continuous or quantized.

- 4. Expressive models require a lot of layers, parameters.

- 5. Lack of flexibility in model design.

- Solution: use approximate inference for intractable models.

- Approximate inference:

- 1. Markov Chain Monte Carlo: generate samples from the exact posterior using a Markov Chain:

- +: very general; exact in the limit of infinite time/computation.

- : computationally expensive; convergence is hard to diagnose.

- 2. Variational inference: approximate the posterior with a tractable distribution; e.g. fully factorized or auto-regressive.

- +: fairly efficient, inference is reduced to optimization w.r.t. dist. parameters.

- : cannot trade computation for accuracy easily.

- Variational inference:

- We approximate the exact posterior  $p_\theta(z|x)$  with a variational posterior  $q_\phi(z|x)$

- We can use any distribution for  $q_\phi(z|x)$  as long as:

- 1. We can sample from it.

- 2. We can compute  $\log q_\phi(z|x)$  and its gradient w.r.t.  $\phi$

Training: The variational posterior induces a variational lower bound  $L_{\theta,\phi}(x)$  on the marginal log-likelihood.

→ we train the model by maximizing  $L_{\theta,\phi}(x)$  w.r.t both

model parameters  $\theta$  and the variational parameter  $\phi$ .

Bounding the marginal log-likelihood:

$$\log p_\theta(x) = \log \int p_\theta(x,z) dz = \log \int q_\phi(z) \frac{p_\theta(x,z)}{q_\phi(z)} dz \geq \int q_\phi(z) \log \frac{p_\theta(x,z)}{q_\phi(z)} dz$$

$$\Rightarrow L_{\theta,\phi}(x) = \mathbb{E}_{q_\phi(z)} [\log \frac{p_\theta(x,z)}{q_\phi(z)}] \leq \log p_\theta(x)$$

- Maximizing likelihood:  $\theta^* = \arg \max \sum_{i=1}^n \log p_\theta(x_i)$

- The gradient for a single point:

$$\nabla_\theta \log p_\theta(x) = \int p_\theta(z|x) \nabla_\theta \log p_\theta(z) dz$$

- Gradient w.r.t. variational parameters:

Two major types of gradient estimator  $\nabla_\phi \mathbb{E}_{q_\phi(z)} [f(z)]$ :

- 1. REINFORCE / likelihood-ratio estimator

- +: very general (discrete & continuous, non-differentiable  $f(z)$ )

- : high variance.

- 2. Reparametrization / path-wise estimator

- +: less general (continuous, differentiable  $f(z)$ )

- +: relatively low variance.
- </div