# **Deep Generative Models**

CS772A: Probabilistic Machine Learning Piyush Rai

## Latent Variable Models for Generation Tasks

• Assume a K-dim latent variable  $\boldsymbol{z}_n$  is transformed to generate to D-dim observation  $\boldsymbol{x}_n$ 



- It is common to use a Gaussian prior for  $\boldsymbol{z}_n$  (though other priors can be used)
- If we use a neural net or GP, such models can generate very high-quality data
  - Take the trained network, generate a random  $m{z}$  from prior, pass it through the model to generate  $m{x}$



Some sample images generated by Vector Quantized Variational Auto-Encoder (VQ-VAE), a state-of-the-art latent variable model for generation



## Factor Analysis and Probabilistic PCA



- FA and PPCA assume f to be a linear model
- $\blacksquare$  In FA/PPCA, latent variables  $\pmb{z}_n \in \mathbb{R}^K$  typically assumed to have a Gaussian prior
  - If we want sparse latent variabled, can use Laplace or spike-and-slab prior on  $z_n$
  - More complex extensions of FA/PPCA use a mixture of Gaussians prior on  $z_n$
- Assumption: Observations  $x_n \in \mathbb{R}^D$  typically assumed to have a Gaussian likelihood
  - Other likelihood models (e.g., exp-family) can also be used if data not real-valued
- Relationship between  $z_n$  and  $x_n$  modeled by a noisy linear mapping

$$\begin{aligned} \mathbf{x}_n &= \mathbf{W} \mathbf{z}_n + \epsilon_n \\ \text{Zero-mean and diagonal or spherical Gaussian noise} \end{aligned} = \sum_{k=1}^{K} \mathbf{w}_k \mathbf{z}_{nk} + \epsilon_n \\ \text{Linear combination of the columns of } \mathbf{W} \end{aligned} \qquad \begin{aligned} p(\mathbf{z}_n) &= \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I}) \\ p(\mathbf{x}_n | \mathbf{z}_n) &= \mathcal{N}(\mathbf{x}_n | \mathbf{W} \mathbf{z}_n, \Psi) \\ \text{Diagonal for FA, spherical for PPCA} \end{aligned}$$
  
**Linear Gaussian Model.**  $\mathbf{W}$ ,  $\mathbf{z}_n$ 's, and  $\Psi$  can be learned (e.g., using EM, VI, MCMC)

# Some Variants of FA/PPCA

Gamma-Poisson latent factor model

Popular for modeling countvalued data (in text analysis, recommender systems, etc) Non-negative priors often give a nice interpretability to such latent variable models (will see some more examples of such models shortly)

- Assumes K-dim non-negative latent variable  $\mathbf{z}_n$  and D-dim count-valued observations  $\mathbf{x}_n$
- ${\mbox{-}}$  An example: Each  $x_n$  is the word-count vector representing a document

 $p(\mathbf{z}_{n}) = \prod_{k=1}^{K} \text{Gamma}(\mathbf{z}_{nk}|\mathbf{a}_{k}, \mathbf{b}_{k}))$  $p(\mathbf{x}_{n}|\mathbf{z}_{n}) = \prod_{d=1}^{D} \text{Poisson}(\mathbf{x}_{nd}|f(\mathbf{w}_{d}, \mathbf{z}_{n}))^{\checkmark}$ 

This is the rate of the Poisson. It should be non-negative,  $\exp(\mathbf{w}_d^{\mathsf{T}} \mathbf{z}_n)$ , or simply  $\mathbf{w}_d^{\mathsf{T}} \mathbf{z}_n$  if  $\mathbf{w}_d$  is also non-negative (e.g., using a gamma/Dirichlet prior on it)

- This can be thought of as a probabilistic non-negative matrix factorization model
- Dirichlet-Multinomial/Multinoulli PCA
  - Assumes K-dim non-negative latent variable  $\mathbf{z}_n$  and D categorical obs  $\mathbf{x}_n = \{\mathbf{x}_{nd}\}_{d=1}^D$
  - An example: Each  $\mathbf{x}_n$  is a document with D words in it (each word is a categorical value)

Also sums to 1

 $p(\mathbf{z}_n) = \text{Dirichlet}(\mathbf{z}_n | \boldsymbol{\alpha})$ 

 $p(\mathbf{x}_{n}|\mathbf{z}_{n}) = \prod_{d=1}^{D} \text{Multinoulli}(\mathbf{x}_{nd}|f(\mathbf{w}_{d},\mathbf{z}_{n}))$ 

This should give the probability vector of the multinoulli over  $x_{nd}$ . It should be non-negative and should sums to 1

VAE\* is a probabilistic extension of autoencoders (AE). An AE is shown below



- The basic difference is that VAE assumes a prior p(z) on the latent code z
  - This enables it to not just compress the data but also generate synthetic data
  - $\hfill\blacksquare$  How: Sample z from a prior and pass it through the decoder
- Thus VAE can learn good latent representation + generate novel synthetic data
- The name has "Variational" in it since it is learned using VI principles

5

# Variational Autoencoder (VAE)



- The Reparametrization Trick is commonly used to optimize the ELBO
- Posterior is inferred only over z, and usually only point estimate on  $\theta$

#### Amortized Inference

- Latent variable models need to infer the posterior  $p(\mathbf{z}_n | \mathbf{x}_n)$  for each observation  $\mathbf{x}_n$
- This can be slow if we have lots of observations because
  - 1. We need to iterate over each  $p(\boldsymbol{z}_n | \boldsymbol{x}_n)$
  - 2. Learning the global parameters needs wait for step 1 to finish for all observations
- One way to address this is via Stochastic VI
- Amortized inference is another appealing alternative (used in VAE and other LVMs too)

 $p(\mathbf{z}_n | \mathbf{x}_n) \approx q(\mathbf{z}_n | \phi_n) = q(\mathbf{z}_n | NN(\mathbf{x}_n; \mathbf{W}))$  If q is Gaussian then the NN will output a mean and a variance

- Thus no need to learn  $\phi_n$ 's (one per data point) but just a single NN with params W
  - This will be our "encoder network" for learning  $\boldsymbol{z}_n$
  - Also very efficient to get  $p(\boldsymbol{z}_*|\boldsymbol{x}_*)$  for a new data point  $\boldsymbol{x}_*$

Variational Autoencoder: The Complete Pipeline

Both probabilistic encoder and decoder learned jointly by maximizing the ELBO





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

# Generative Adversarial Network (GAN)

- GAN is an implicit generative latent variable model
- Can generate from it but can't compute p(x) the model doesn't define it explicitly
- GAN is trained using an adversarial way (Goodfellow et al, 2013)



Thus can't train

Unlike VAE, no explicit parametric

likelihood model p(x|z)

Generative Adversarial Network (GAN)

The GAN training criterion was

 $\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{data}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$ 

- With G fixed, the optimal D (exercise)  $D_{G}^{*}(x) = \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)}$ Distribution of synthetic data
- Given the optimal D, The optimal generator G is found by minimizing

$$V(D_{g}^{*},G) = \mathbb{E}_{x \sim p_{data}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{g}(x)}{p_{data}(x) + p_{g}(x)} \right]$$

$$Jensen-Shannon$$

$$divergence between$$

$$p_{data} \text{ and } p_{g}.$$

$$Minimized \text{ when}$$

$$p_{g} = p_{data}$$

$$P_{g} = p_{data}$$

$$P_{g} = p_{data}$$

$$J(x) = \mathbb{E}_{x \sim p_{data}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{g}(x)}{p_{data}(x) + p_{g}(x)} \right] - \log 4$$

$$U(x) = \mathbb{E}_{x \sim p_{data}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{g}(x)}{p_{data}(x) + p_{g}(x)} \right] - \log 4$$

$$U(x) = \mathbb{E}_{x \sim p_{data}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{g}(x)}{p_{data}(x) + p_{g}(x)} \right] - \log 4$$

$$U(x) = \mathbb{E}_{x \sim p_{data}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{data}(x)}{p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{g}(x)}{p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[ \log \frac{p_{g}(x)}{p_{g}(x)} \right]$$

#### GAN Optimization $\min_{D} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{data}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z}))]$ The GAN training procedure can be summarized as $\theta_a$ and $\theta_d$ denote the params of the deep neural nets $\prec$ defining the generator and discriminator, respectively **1** Initialize $\boldsymbol{\theta}_a, \boldsymbol{\theta}_d$ ; 2 for each training iteration do In practice, for stable training, we run K > 1 steps of for K steps do \_\_\_\_\_\_ optimizing w.r.t. D and 1 step of optimizing w.r.t. G 3 Sample minibatch of M noise vectors $\mathbf{z}_m \sim q_z(\mathbf{z})$ ; 4 Sample minibatch of M examples $\mathbf{x}_m \sim p_D$ ; 5 Update the discriminator by performing stochastic gradient *ascent* using this gradient: 6 $\nabla_{\boldsymbol{\theta}_d} \frac{1}{M} \sum_{m=1}^{M} \left[ \log D(\mathbf{x}_m) + \log(1 - D(G(\mathbf{z}_m))) \right]. ;$

- 7 Sample minibatch of M noise vectors  $\mathbf{z}_m \sim q_z(\mathbf{z})$ ;
- 8 Update the generator by performing stochastic gradient *descent* using this gradient:  $\nabla \theta_g \frac{1}{M} \sum_{m=1}^{M} \log(1 - D(G(\mathbf{z}_m))). ;$

9 Return  $\boldsymbol{\theta}_g, \, \boldsymbol{\theta}_d$ 

In practice, in this step, instead of minimizing

 $\log(1 - D(G(z)))$ , we maximize  $\log(D(G(z)))$ 

Reason: Generator is bad initially so discriminator will always predict correctly initially and log(1 - D(G(z))) will saturate

#### GANs that also learn latent representations

12

- The standard GAN can only generate data. Can't learn the latent z from x
- Bidirectional GAN\* (BiGAN) is a GAN variant that allows this



# **Evaluating GANs**

- Two measures that are commonly used to evaluate GANs
  - Inception score (IS): Evaluates the distribution of generated data
  - Frechet inception distance (FID): Compared the distribution of real data and generated data
- Inception Score defined as  $\exp(\mathbb{E}_{x \sim p_q}[\mathrm{KL}(p(y|x)||p(y))])$  will be high if
  - Very few high-probability classes in each sample x: Low entropy for p(y|x)
  - We have diverse classes across samples: Marginal p(y) is close to uniform (high entropy)
- FID uses extracted features (using a deep neural net) of real and generated data
  - Usually from the layers closer to the output layer
- These features are used to estimate two Gaussian distributions

Using real data  $\mathcal{N}(\mu_R, \Sigma_R)$   $\mathcal{N}(\mu_G, \Sigma_G)$  Using generated data

• FID is then defined as FID =  $|\mu_G - \mu_R|^2 + \text{trace}(\Sigma_G + \Sigma_R - (\Sigma_G \Sigma_R)^{1/2})$ 

These measures can also be used for evaluating other deep gen models like VAEML

High IS and low FID is desirable

Both IS and FID measure how

realistic the generated data is

## GAN: Some Issues/Comments

- GAN training can be hard and the basic GAN suffers from several issues
- Instability of training procedure
- Mode Collapse problem: Lack of diversity in generated samples
  - Generator may find some data that can easily fool the discriminator
  - It will stuck at that mode of the data distribution and keep generating data like that



GAN 1: No mode collapse (all 10 modes captured in generation)

GAN 2: Mode collapse (stuck on one of the modes)

Some work on addressing these issues (e.g., Wasserstein GAN, Least Squares GAN, etc)