

Deep Generative Architectures

Oliver Schulte School of Computing Science Simon Fraser University Introduction to Deep Learning

Generative Probabilistic Models

- Supervised learning models P(y | x): the probability of one or more target variables y given input variables x
- Generative model P(x): just model the distribution of the inputs

Example: synthetic outputs

- See <u>tutorial</u>
- And <u>face generator</u>



Figure 7.10: VAE Mnist digits generated from scratch

How can we get an NN to generate output? Doesn't it just map inputs to outputs?

General Neural Network Approach

- Input random vector z to neural net
 - Typically from a Gaussian bell curve distribution
- Network maps **z** to output vector **x**
- Intuitively, the output should be like the observations x



p(Z) = Gaussian bell curve

Simple Example Mappings

Latent distribution	+1 distribution	
Uniform over [0,1]	Uniform over [1,2]	
Gaussian with mean 0 and standard deviation 1	standard deviation 1 Gaussian with mean 1	
	and standard deviation 1	



Generate Ring from Gaussian



Samples from Gaussian distribution of 2D **Z**

Samples from output distribution

How to Train Deep Generative Models?

- What training objective?
- 2 main ideas
 - Variational Auto-Encoder
 - Generative Adversarial Network

Variational Auto-Encoder (VAE)

VAE Objective Approximates Log-Likelihood

- Loss function for input $\mathbf{x} = -\ln(P(\mathbf{x}))$
 - The negative log-likelihood is the standard loss for generative models
- In the decoder architecture $P(\mathbf{x}) \approx \int_{\mathbf{z}} 1(f(\mathbf{z}) = \mathbf{x}) p(\mathbf{z}) d\mathbf{z}$
 - 1(f(z) = x) returns 1 if the decoder maps random to x, 0 o.w.
- Integral is intractable
- VAE architecture is designed so that <u>training loss</u> <u>approximates integral negative log-likelihood</u>



p(Z) = Gaussian bell curve

VAE Architecture

- Combines Auto-Encoding with Generative Modeling
- High-level idea
- Training: Design a non-deterministic version of an auto-encoder
 - Can produce multiple outputs **x**-out for the same input **x**-in
- 2. Testing: Generate <u>random input</u> **z**, produce (non-deterministic) output





Example: Non-deterministic output



original

reconstruction

VAE Training Architecture







Example Application

Embeddings for Hockey Players

Contextualized Player Embedding

- The set-up: we see a hockey game sequence.
- Using a VAE, compute a contextualized player embedding as a function of the game sequence.
- The VAE generates a distribution $P(player_t = i | sequence up to time t)$



Combining LSTM with VAE

- Sequence is processed using LSTM
- VAE generates probability distribution over players



Accuracy in Player Recognition

Given a match sequence, can you predict which player is likely to have the puck now?

Method	Accuracy	Log-likelihood
LSTM	12.41%	-3.131
LSTM + VAE	48%	-2.228

- The player embeddings help with applications like predicting
 - Whether a shot will lead to a goal
 - The final game outcome

The Evidence Lower Bound

The Mathematics Behind VAEs

Background

- Monte Carlo Sampling
- The KL Divergence

Sampling Example

- Area of unit circle = π
- Estimate π by randomly sampling points in square, checking if they lie within circle. <u>Demo</u>



Sample Complexity

• Basic Problem: Sampling can require <u>many data points</u> for accurate estimate



Kulback-Leibler Divergence

- KLD: $D(q | |p) = \sum_{j} q(\mathbf{x}_{j}) \log(q(\mathbf{x}_{j})/p(\mathbf{x}_{j}))$ = $E_{q}[\log(q(\mathbf{x})/p(\mathbf{x}))]$
- Entropy: $H(q) = -\sum_{j} q(\mathbf{x}_{j}) \log(q(\mathbf{x}_{j}))$ = $-E_{q}[\log(q(\mathbf{x})]$
- Cross-Entropy: $H(q,p) = -\sum_{j} q(\mathbf{x}_{j}) \log(p(\mathbf{x}_{j}))$ = $-E_{q}[\log(p(\mathbf{x})]$
- Exercise: Prove that D(p | |q) = H(q,p) H(q)
- Exercise: Recall the "cross-entropy" of Assignment $\sum_{j=1}^{N} y_j ln(p_j^+) + (1 - y_j)ln(1 - p_j^+)$ What is the relationship between this formula and H(q,p)?

Likelihood: The True Objective

- Recall that the standard objective for generative modelling is the data likelihood
 Π_j P(**x**_j; f) = ∫_{**z**} 1(f(**z**) = **x**_j) p(**z**) d**z** for observed data points **x**₁,...,**x**_j
- We could estimate this integral for each data point using <u>Monte Carlo estimation.</u>
 - 1. Sample $k \mathbf{z}$ points,
 - 2. apply f
 - 3. Calculate the mean estimate: average $f(\mathbf{z}_1), \ldots, f(\mathbf{z}_k)$

Sampling With Functions



Samples from Gaussian distribution of 2D P(**Z**)

Samples from output distribution P(**X**)

- If f is invertible, then for each output **x**, there is a unique $f^{-1}(z)$ that generates it.
- Even if f is not invertible, typically for each output **x**, few samples will generate **x**. **Generative Models**

Invertible Function Case

- We can replace $P(\mathbf{x}; f) = \int_{\mathbf{z}} 1(f(\mathbf{z}) = \mathbf{x}) p(\mathbf{z}) d\mathbf{z} = \int_{\mathbf{z}} p(f^{-1}(\mathbf{x})) d\mathbf{z}$
- I.e., just consider probability of the random input that generates the observed output **x**
- How to compute $f^{-1}(\mathbf{x})$?
 - Learn another neural network for f⁻¹!



Probabilistic Version

- Even if the generating function *f* is not invertible, few samples z have a significant likelihood of generating an observed data point x.
- We can measure this with the *posterior probability* p(z | x): given that x was observed, what is the probability that z generated x.
- The posterior probability is the inverse of a *conditional* probability model $p(\mathbf{x} | \mathbf{z})$.
- $P(\mathbf{x}) = \int_{\mathbf{z}} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z} \approx \int_{\mathbf{z}} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z} | \mathbf{x}) d(\mathbf{z})$



Approximate Posterior

- Given the approximation $P(\mathbf{x}) \approx \int_{\mathbf{z}} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z} | \mathbf{x}) d(\mathbf{z})$
- How can we compute $p(\mathbf{z} \mid \mathbf{x})$?
- Learn a neural network to output (the parameters of) $p(\mathbf{z} \mid \mathbf{x})!$
- Hard to get this exactly so we try to learn an *approximate* posterior q(z | x) ≈ p(z | x)



Monte Carlo Estimation of Data Likelihood

• $P(\mathbf{x}) \approx \int_{\mathbf{z}} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z} | \mathbf{x}) d(\mathbf{z})$

 $\boldsymbol{\approx} \operatorname{E}_{\mathbf{z} \thicksim q(\mathbf{z} \mid \mathbf{x})} p(\mathbf{x} \mid \mathbf{z})$

- For each data point **x**
 - 1. Apply the encoder to find an <u>approximate posterior</u> distribution $q(\mathbf{z} \mid \mathbf{x}_j)$
 - 2. Sample \mathbf{z} from $q(\mathbf{z} | \mathbf{x}_j)$
 - 3. Apply the decoder to compute $p(\mathbf{x} | \mathbf{z})$
 - 4. Average the estimates from the z samples to get an estimated P(x)
- See <u>VAE tutorial</u>
- What is the relationship between the encoder-decoder approximations and the true data likelihood P(**x**)?

The Evidence Lower Bound

The following relationship holds for *any* approximate posterior q(z | x), prior p(z), conditional probability p(x | z), unconditional probability p(x)

 $\log(\mathbf{p}(\mathbf{x})) \ge E_{\mathbf{z} \sim q(\mathbf{z} \mid \mathbf{x})} [\log(\mathbf{p}(\mathbf{x} \mid \mathbf{z}))] - KLD(q(\mathbf{z} \mid \mathbf{x}) \mid |\mathbf{p}(\mathbf{z}))$

True log- prob	Sample z from approximate posterior (encoder)	Compute log- probability of x given z (decoder)	Penalize differences between app. Posterior and true prior
-------------------	---	--	--

- The bigger the RHS, the better our approximation.
- So we try to maximize it: argmax p,q $\sum_{j} E_{\mathbf{z} \sim q(\mathbf{z} \mid \mathbf{x}_j)} [log(p(\mathbf{x}_j \mid \mathbf{z}))] - KLD(q(\mathbf{z} \mid \mathbf{x}_j) \mid | p(\mathbf{z}))$

Implementing the ELBO objective

- The ELBO is a very general result, known for decades.
- The VAE is a neural architecture proposed by Kingma and Welling 2014 that fills in q and p as follows.
- $q(\mathbf{z} | \mathbf{x}) = Normal(\mu(\mathbf{x}), \sigma^2(\mathbf{x}))$
 - $\bullet~\mu$ and σ are vectors of the same dimension as z
 - They are computed by a neural network (encoder)
- $p(\mathbf{x} | \mathbf{z}) = \text{Normal}(f(\mathbf{z}), \sigma^2 * I)$
 - f(z) is a vector of the same dimension as x
 - σ^2 is a hyper-parameter
 - f(z) is computed by a neural network (decoder)
- <u>Visualization</u>

KLD for Gaussians

- The ELBO requires us to evaluate $KLD(q(\mathbf{z} | \mathbf{x}_j) | | p(\mathbf{z}))$
- For two multi-variate Gaussians $N_Q(m_Q, \Sigma_Q)$ and $N_p(m_p, \Sigma_p)$, the KLD divergence is given by the formula $D(N_Q | | N_p) = \frac{1}{2}$ $\{tr (\Sigma^{-1}_p \Sigma_q) + ln | \Sigma_p | - ln | \Sigma_Q | + (m_p - m_Q)^T \Sigma^{-1}_p (m_p - m_Q)^T \} - d/2$ where
- d is the latent embedding dimension
- tr is the matrix trace (sum of diagonal elements)
- In a VAE, we have $N_p = (0, I)$ and Σ_q is diagonal with entries σ_i^2 for i = 1, ..., d
- Exercise: Show that $D(N_Q | | N_p) = \frac{1}{2} \{ \sum_i \sigma_i^2 - \ln(\sigma_i^2) + m_i^2 \} - d/2 \text{ where } m = m_Q$

Generative Adversarial Models

GANs

GAN Architecture

- Recall the basic generative NN architecture
- How to train?
- VAE: approximate log-likelihood of the observations
- GAN: train the generator so that synthetic generated examples cannot be distinguished from actual observations





Intuition: Taste Test

- Coke is the real thing
- Pepsi wants to imitate coke
- How does Pepsi know they have succeeded?
- When a blind taste test cannot tell the difference!





Chatbots and the Turing Test

Alan Turing (1950) "Computing Machines and Intelligence"

- How can we know if a machine is intelligent?
- Test if it can fool a human!
- <u>Loebner Prize</u>
- ChatGPT

Distinguishing Uniform and Gaussians

Fake data: uniform [-8,8]
 Real data: Gaussian with mean 5, standard deviation 1



The Tester

- How can we test whether a network generates realistic outputs?
- Train another classifier network to distinguish synthetic from real examples!
- Discriminator *D* outputs $o_D(x) = P(x \text{ is real example})$
- Why can we not just backpropagate to train the generator?



The Loss Functions

- Discriminator objective function given real examples *R* and fake examples *F* $V(D,G) = mean_r ln(o_D(r)) + mean_f ln(1-o_D(r))$
- The discriminator wants to maximize V, the generator to minimize V
- Minmax problem: GANs are tricky to train

	Discriminator			Generator	
x	Output	Objective	Loss	Objective	Loss
Real r	$o_D(r)$	$\ln(o_D(\mathbf{r}))$	$-\ln(o_D(\mathbf{r}))$	$-\ln(o_D(\mathbf{r}))$	$\ln(o_D(\mathbf{r}))$
Fake f	$o_D(f)$	$\ln(1-o_D(f))$	$-\ln(1 - o_D(f))$	$-\ln(1 - o_D(f))$	$\ln(1-o_D(f))$

Other Loss Functions

- Other loss functions can be considered to help with training, e.g. in the book:
- The generator wants $o_D(f)$ to be big
 - So maximize $\ln o_D(f)$ or minimize $-\ln o_D(f)$
- However it is <u>more usual</u> to ensure the zero-sum condition: loss(discriminator) = - loss(generator)

Book	Discriminator		Generator		
x	Output	Objective	Loss	Objective	Loss
Real r	$o_D(r)$	$\ln(o_D(\mathbf{r}))$	$-\ln(o_D(\mathbf{r}))$		
Fake f	$o_D(f)$	$\ln(1-o_D(f))$	$-\ln(1 - o_D(f))$	$\ln(o_D(f))$	$-\ln(o_D(f))$

Examples and Demos

- <u>Style GAN</u>
- <u>Deep Fakes</u>
- <u>GANlab</u>
- <u>nvidia-ai-playground</u>

Conclusion

- Generative models generate outputs without an input
- Basic idea: map a random input to an output
- Without a target output, what is the training objective?
- Variational AE: approximate the log-likelihood of the observed data **x**
 - \bullet Also produces an embedding of the input ${\bf x}$
- Generative Adversarial Network: generate synthetic outpus that cannot be distinguished from actual observations by an adversarial classifier network