

F23-CS-236-01 Midterm

Jason Alexander Chan

TOTAL POINTS

59.4 / 95

QUESTION 1

Question 1 20 pts

1.1 a 0 / 2

- 0 pts Correct

- 1 pts Incorrect or insufficient explanation.

✓ - 2 pts *Incorrect or missing response.*

✓ - 0 pts *Correct*

- 1 pts Incorrect or no explanation

- 2 pts Incorrect

- 1 pts Correct idea but the answer should be false

1.2 b 2 / 2

✓ - 0 pts *Correct*

- 1 pts Incorrect or insufficient explanation.

- 2 pts Incorrect or missing response.

1.7 g 2 / 2

✓ - 0 pts *Correct*

- 2 pts Click here to replace this description.

1.3 c 2 / 2

✓ - 0 pts *Correct*

- 1 pts Incorrect or insufficient explanation.

- 2 pts Incorrect or missing response.

1.8 h 0 / 2

- 0 pts *Correct*

✓ - 2 pts *Click here to replace this description.*

- 1 pts incorrect or no explanation

1.4 d 2 / 2

✓ - 0 pts *Correct*

- 1 pts Incorrect or insufficient explanation.

- 2 pts Incorrect or missing response.

1.9 i 1 / 2

- 0 pts *Correct*

- 2 pts Click here to replace this description.

✓ - 1 pts *incorrect or missing explanation*

1.10 j 2 / 2

✓ - 0 pts *Correct*

- 2 pts Click here to replace this description.

1.5 e 0 / 2

- 0 pts *Correct*

- 1 pts Incorrect or insufficient explanation.

✓ - 2 pts *Incorrect or missing response.*

QUESTION 2

Question 2 10 pts

2.1 a 1.6 / 2

✓ + 0.4 pts AR

1.6 f 2 / 2

✓ + 0.4 pts VAE
✓ + 0.4 pts EBM
+ 0.4 pts no GAN
✓ + 0.4 pts no Flow
+ 0 pts Click here to replace this description.

✓ + 0.4 pts VAE
✓ + 0.4 pts Flow
✓ + 0.4 pts GAN
✓ + 0.4 pts EBM
+ 0 pts Incorrect

2.2 b 2 / 2

✓ + 0.4 pts AR
✓ + 0.4 pts Flow
✓ + 0.4 pts no EBM
✓ + 0.4 pts no VAE
✓ + 0.4 pts no GAN
+ 0 pts Click here to replace this description.

QUESTION 3

Question 3 18 pts

2.3 C 2 / 2

✓ + 0.4 pts AR
✓ + 0.4 pts VAE
✓ + 0.4 pts Flow
✓ + 0.4 pts GAN
✓ + 0.4 pts EBM
+ 0 pts Click here to replace this description.

3.1 a 0 / 3

- 0 pts Correct
✓ - 3 pts Wrong
- 1 pts Correct but fails to mention prior needs a tractable likelihood

2.4 d 0.8 / 2

+ 2 pts Correct
✓ + 0.4 pts AR
+ 0.4 pts VAE
+ 0.4 pts Flow
+ 0.4 pts GAN
✓ + 0.4 pts EBM
+ 0 pts Incorrect

3.3 C 0 / 3

- 0 pts Correct
✓ - 3 pts incorrect
- 3 pts missing

2.5 e 2 / 2

+ 2 pts Correct
✓ + 0.4 pts No AR

3.4 d 3 / 3

✓ - 0 pts Correct
- 3 pts reasoning incorrect/missing
- 2 pts correct idea, not rigorous

3.5 e 3 / 3

✓ - 0 pts Correct
- 3 pts Incorrect
- 2 pts Does not include an explanation.
- 1 pts No mentioning of reparametrization trick

on Z / Z is still a gaussian distribution.

3.6 f 0 / 3

- 0 pts Correct

✓ - 3 pts Incorrect

- 1 pts No explanation.

- 1 pts Wrong explanation.

- 0.5 pts Minor error(i.e. x,y still in the equation/did not specify variables in P_xy/minor error in calculation)

- 4 pts Missing

- 0 pts Correct

QUESTION 4

Question 4 10 pts

4.1 a 2 / 2

✓ - 0 pts Correct: $1/3 * 1/5 = 1/15$

- 0.5 pts Partially wrong: give the correct expression of x but not the answer 1/15

- 1 pts Partially wrong: either Jacobian (1/3) or $\$\$p_z(f^{-1}(x))\$$(1/5)$ is wrong

- 2 pts Wrong

QUESTION 5

Question 5 11 pts

5.1 a 0 / 2

- 0 pts Correct

- 1 pts Partially correct

✓ - 2 pts Incorrect/missing

5.2 b 1 / 2

- 0 pts Correct

✓ - 1 pts Partially Correct

- 2 pts Incorrect/missing

5.3 C 2 / 2

✓ - 0 pts Correct

- 2 pts Incorrect/missing

- 1 pts Partially Correct

5.4 d 4 / 4

✓ - 0 pts Correct

- 3 pts Missing expectation for the KL-Divergence

- 3 pts Missing / unclear dependencies on t

- 2 pts Incomplete

- 4 pts Answer missing

4.3 C 2 / 2

✓ - 0 pts Correct: Yes, because there is a unique inverse

- 1 pts True/False is correct but the explanation is wrong or missing

- 2 pts Wrong

4.4 d 0 / 4

✓ - 2 pts Incorrect Jacobian

✓ - 2 pts Incorrect joint pdf

5.5 e 1 / 1

✓ - 0 pts Correct

- 1 pts Incorrect. Correct answer: "Yes, the

sampling step can be rewritten as a deterministic transformation of standard Gaussians by unrolling the RNN."

QUESTION 6

Question 6 10 pts

6.1 a 1 / 2

- 0 pts Correct
- ✓ - 1 pts Incorrect logic / missing logic / missing mathematical formalization
- 2 pts Incorrect / missing

6.2 b 2 / 2

- ✓ - 0 pts Correct
- 1 pts Minor Error / Did not simplify / Did not follow instructions
- 2 pts Incorrect

6.3 c 2 / 2

- ✓ - 0 pts Correct
- 1 pts Minor Error / Did not simplify / Did not follow instructions
- 2 pts Incorrect

6.4 d 1 / 4

- 0 pts Correct
- ✓ - 1 pts *D(G(z)) value at the beginning of training*
- 2 pts magnitudes of gradients in training
- ✓ - 1 pts *training speed*
- 4 pts Incorrect or missing
- ✓ - 1 pts *Minor Error*

QUESTION 7

Question 7 16 pts

7.1 a 2 / 2

- ✓ - 0 pts Correct
- 1 pts on the right track, but did not explicitly say $Z(1)=1$, or wrong $Z(1)$
- 1 pts asking for $Z(1)$ instead of $\log Z(1)$
- 2 pts incorrect

7.2 b 0 / 2

- 0 pts Correct
- ✓ - 2 pts incorrect
- 1 pts right track, but didn't answer explicitly tractable or not

7.3 c 2 / 2

- ✓ - 0 pts Correct
- 2 pts incorrect
- 1 pts too vague, not explicitly $0 < T < 1$
- 1 pts T should be strictly less than 1

7.4 d 1 / 2

- 0 pts Correct
- ✓ - 1 pts *Sampling from p(x) is tractable*
- 1 pts Sampling from $p^T(x)$ is intractable

7.5 e 2 / 2

- ✓ - 0 pts Correct
- 1 pts Minor incorrect
- 2 pts Incorrect

7.6 f 0 / 4

- 0 pts Correct
- 2 pts Incomplete/slightly incorrect
- ✓ - 4 pts *Incorrect/Not answered*
- 1 pts Minor missing steps

7.7 g 2 / 2

✓ - 0 pts Correct

- 2 pts incorrect

Name: Jason Alayandar SUNet ID: jchana7 Stanford.edu

Note: Partial credit will be given for partially correct answers. Zero points will be given to answers left blank.

Total score: / 96 $\frac{180 \text{ marks}}{180 \text{ marks}} = 1.0 \text{ marks/10 pts}$

Question	Score	Question	Score
10.40am 1	/ 20	12.20 5	/ 11
11 am 2	/ 10	12.40 6	/ 10
11.40am 3	/ 16	1.00pm 7	/ 10
12pm 4	/ 10		

- While the faculty alone has the right and obligation to set academic requirements, the students and faculty will work together to establish optimal conditions for honorable academic work.

- The faculty on its part maintains its confidence in the honor of its students by refraining from profiting from violations and from taking unusual and unreasonable precautions to prevent the forms of dishonesty mentioned above. The faculty will also avoid, as far as practicable, academic procedures that create temptations to violate the Honor Code.

- They will do their share and take an active part in seeing to it that others as well as themselves uphold the spirit and letter of the Honor Code.
- That they will not give or receive aid in examinations; that they will not give or receive unpermitted aids in class work, in the preparation of reports, or in any other work that is to be used by the instructor as the basis of grading;
- That they will not give or receive aid in examinations; that they will not give or receive unpermitted aids in class work, in the preparation of reports, or in any other work that is to be used by the instructor as the basis of grading;
- That they will not give or receive aid in examinations; that they will not give or receive unpermitted aids in class work, in the preparation of reports, or in any other work that is to be used by the instructor as the basis of grading;

- The Honor Code is an undertaking of the students, individually and collectively:

The Honor Code is the University's statement on academic integrity written by students in 1921. It articulates University expectations of students and faculty in establishing and maintaining the highest standards in academic work:

Stanford University Honor Code

This exam is worth 95 points. You have 3 hours to complete and submit it. You are allowed to consult notes, books and use a laptop. But no communication or network access is allowed. Use of Large Language Models (LLMs) is also not allowed. Good luck!

CS 236, Fall 2023 Midterm Exam

I attest that I have not given or received aid in this examination, and that I have done my share and taken an active part in seeing to it that others as well as myself uphold the spirit and letter of the Honor Code.

Signature:



1. [20 points total] True/False

For each of the statements below, slate True or False. Explain your answer for full points.

(a) [2 points] Autoregressive models:

affect the initial number of parameters required to represent exactly the joint distribution.

For discrete autoregressive models, the order of variables in the chain rule factorization cannot

True. The number of terms to represent the exact joint

autoregressive Model: Let $p(x)$ be an arbitrary autoregressive model over a set of discrete random variables. The probability assigned by p to a sequence $x = (x_1, x_2, \dots, x_n)$ of length n is always at least as large as the probability assigned by p to any sequence $(x_1, x_2, \dots, x_n, x_{n+1})$ of length $n + 1$ of which x is a prefix.

(b) [2 points] Autoregressive Models: Let $p(x)$ be an arbitrary autoregressive model over a set of discrete random variables.

length n is always at least as large as the probability assigned by p to a sequence $(x_1, x_2, \dots, x_n, x_{n+1})$ of length $n + 1$ of which x is a prefix.

(c) [2 points] VAE: In a VAE with two binary latent variables $z_1 \in \{0, 1\}; z_2 \in \{0, 1\}$ and a Gaussian decoder $p(x|z_1, z_2)$, the marginal likelihood $p(x)$ is intractable to compute.

(d) [2 points] ELBO: Let the inference distribution of a VAE be $q_\phi(z|x)$ and prior distribution $p(x)$, the gap between the MLE objective and the ELBO lower bound for a datapoint x is exactly $D_{KL}(q_\phi(z|x) \| p(z))$.

(e) [2 points] ELBO: Alice and Bob train two VAEs on the same dataset. The VAEs are identical, except for the priors over the latents which are different. Alice's model also achieves higher log-likelihood on the ELBO on the training set. This implies Alice's model achieves a better (higher) ELBO than or equal to Bob's.

(f) [2 points] Normalizing-Flow: A normalizing flow based model with latent variables $z \in \mathbb{R}^m$ can be used to model with tractable exact likelihood evaluation) any random vector x of dimension less than or equal to 100.

(g) [2 points] GAN: A VAE with a Gaussian prior and a Gaussian decoder can be trained as the generator in a GAN (discarding the encoder and using a separate classifier as a discriminator).

(h) [2 points] GAN: A conditional GAN can be implemented by adding the class information as an additional input to the Generator $G(z)$, without changes to the Discriminator $D(x)$.

True. GAN can model the distribution to learn the condition.

False. GAN is negative. GANs with loss early

(i) [2 points] EBM: For general energy-based models, the energy $E(x)$ must be positive for every x .

- (j) [2 points] EBM: There exist energy-based models that can be equivalently described as an autoregressive model, regardless of the variable ordering and parameterization of the conditionals.
- Take: Autoregressive models can be described as energy models.
because they are products of reward-based objects.
Then they can be described as an EBM.
Therefore can evaluate as an EBM.

2. [10 points total] Comparison of Models
- (a) [2 Points] Which of these models can be used to model both discrete and continuous data?
- So far we discussed five major types of generative models: autoregressive models (AR), variational autoencoders (VAE), flow models (Flow), generative adversarial networks (GAN), and energy-based models (EBM). In the following questions, we will compare their strengths and weaknesses. Brief justifications (1 sentence) are sufficient for full points.
- (b) [2 Points] Suppose we are interested in exactly evaluating the likelihood of a data point under the trained model. In which of these models can we exactly evaluate a data point's likelihood in an efficient way? (i.e., in a time polynomial in the number of dimensions of a sample, such as in linear time)?
- (c) [2 Points] Why do these models can be used to model conditional distributions (e.g., trained to map captions to images given a suitable paired dataset)?
- (d) [2 Points] Which of these models can be used to solve "in-painting" tasks (i.e., fill in missing values in a datapoint such as completing a sentence or an image where some parts have been masked out), assuming access to infinite compute?
- (e) [2 Points] Suppose we are interested in learning Distributional Models that e.g. characterize, EBM's have a unique/in-painting feature from eg. EBM's, specifically distribution models can generate samples from test data. AR, Flow, EBM, GAN, EBM for example can map zero-zero features to images given a suitable paired dataset?
- VAE, GAN and EBM require sampling training data to approximate likelihood
- AR, Flow
- VAE, GAN, EBM
- EBMs, specifically distribution models can generate samples from test data.
- AR, Flow, EBM, GAN for example can map zero-zero features to images given a suitable paired dataset?
- which of these models are most appropriate for this task, and why?
- AR could because it's problem formulation includes modeling the world. VAE, EBM with modification, AR with modification, GAN with modification, EBM with modification, AR, Flow, VAE, EBM

(e) [3 points] Let $\mathbf{x} \in \mathbb{R}^D$ denote the inputs, \mathbf{z} the latent variables, $p_\theta(\mathbf{x}|\mathbf{z})$ the generative model, $N(\mathbf{z}; \mu_\phi(\mathbf{x}), \text{diag}(q_\phi(\mathbf{x})^2))$. Consider an alternative inference model representing a Gaussian distribution $p(\mathbf{z})$ the prior and $q_\phi(\mathbf{z}|\mathbf{x})$ as the basic inference model. Formulate the ELBO formula such that $q_\phi(\mathbf{z}|\mathbf{x})$ is at least as good as $q_\theta(\mathbf{z}|\mathbf{x})$.

Many A could result as I if a deep layer. Here the ELBO with matrix A adds more modeling flexibility to learning. The VAE

$$\text{ELBO}(\mathbf{x}; p, q) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})]$$

and A is optimized over the set of all $D \times D$ matrices.

where

$$\max_{\theta, \phi, A} \text{ELBO}(\mathbf{x}; p_\theta, q_\phi, A) \geq \max_{\theta, \phi} \text{ELBO}(\mathbf{x}; p_\theta, q_\phi)$$

Show that the best evidence lower bound we can achieve with $q_\phi(\mathbf{z})$ is at least as tight as the best one we can achieve with $q_\theta(\mathbf{z})$.

$N(\mathbf{z}; \mu_\phi(\mathbf{x}), \text{diag}(q_\phi(\mathbf{x})^2))$. Consider an alternative inference model $q_\phi(\mathbf{z}|\mathbf{x})$ also representing a Gaussian distribution with parameters $N(\mathbf{z}; \mu_\theta(\mathbf{x}), \text{diag}(q_\theta(\mathbf{x})^2))$ where A is a $D \times D$ matrix.

(d) [3 points] Let $\mathbf{x} \in \mathbb{R}^D$ denote the inputs, \mathbf{z} the latent variables, $p_\theta(\mathbf{x}|\mathbf{z})$ the generative model,

VAEs a measure of the same model, even the negative ELBO

is not possible for both to claim global optimum of the ELBO objective?

both claim to have obtained a global optimum of the ELBO objective. When they evaluate their models on a less-discriminative encoder map \mathbf{x} to different distributions over the latents. Is this possible if they truly have both obtained a global optimum of the ELBO objective?

(c) [3 points] Alice and Bob are training two identical VAEs on the same data distribution. They

therefore draw no difference in reconstruction X when the VAE

Yes because the upcode is perfect compression

to the decoder. Is \mathbf{x} distributed according to $p_\theta(\mathbf{x}|\mathbf{z})$ by feeding \mathbf{z} to the encoder, then draw a sample $\mathbf{x}' \sim p_\theta(\mathbf{x}|\mathbf{z}')$ from $p_\theta(\mathbf{x}|\mathbf{z}')$ and $p_\theta(\mathbf{x})$ is zero. Suppose we start from a real datapoint \mathbf{x} sampled from $p_\theta(\mathbf{x}|\mathbf{z})$ between $p_\theta(\mathbf{x}|\mathbf{z})$ and $p_\theta(\mathbf{x})$ is zero. Note this implies that the KL divergence

optimization, arbitrary flexible decoder and encoder). Note this implies that the KL divergence between $p_\theta(\mathbf{x}|\mathbf{z})$ and $p_\theta(\mathbf{x})$ is zero. Suppose we start from $p_\theta(\mathbf{x}|\mathbf{z})$, feed it to the encoder

and obtain a global optimum of the ELBO parameterized by θ for encoder and θ for decoder.

(b) [3 points] Suppose we have trained a VAE parameterized by θ for encoder and θ for decoder.

real world in practice know GAN having style can be unifiable.

however the VAE optimizing ELBO, however this would be

yes if GAN learning object faultless back propagation when

can we train this VAE by optimizing the ELBO?

Adversarial Network instead of standard Gaussian (i.e., samples from the prior are produced by a trainable generator G_θ), and where the encoder $q_\phi(\mathbf{z}|\mathbf{x})$ and decoder $p_\theta(\mathbf{x}|\mathbf{z})$ are Gaussian.

(a) [3 points] Suppose we are training a VAE where the prior $p_\theta(\mathbf{z})$ is parameterized by a Generative

Autoencoder model where as usual \mathbf{x} denotes observed variables and \mathbf{z} denotes latent variables. For

each of the following questions briefly explain your answer for full points.

3. [18 points total] Variational Autoencoders Basics In this question, we will consider a Variational

- So \mathbf{z} is generated by Gaussian \mathbf{w}^ϕ . Using the performance metric which we would compute $\mathbf{z} = \mathbf{B}(\mathbf{w}^\phi + \mathbf{e}) = \mathbf{B}\mathbf{w}^\phi + \mathbf{B}\mathbf{e}$
 because \mathbf{B} is a linear transformation
- (f) [3 points] Compare $\max_{\theta, \phi, A} \text{ELBO}(\mathbf{x}; p_\theta, q_\phi, A)$ to $\max_{\theta, \phi, B} \text{ELBO}(\mathbf{x}; p_\theta, q_\phi, B)$. Is one always larger than the other without assumptions on p_θ and q_ϕ ?
- linear transformation, one to two form a single training set
 If some \mathbf{y} would be the same. They both introduce a

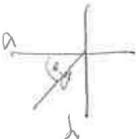
$$\begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} R \cos \theta \\ R \sin \theta \end{bmatrix}$$

general form formula

$$P_{\theta}(x; \theta) = P_x(f_\theta(x)) \quad \text{for } \frac{\partial}{\partial \theta} \cos(\theta) = -\sin(\theta), \quad \frac{\partial}{\partial \theta} \sin(\theta) = \cos(\theta)$$

*Hint: Potentially useful formula:
of $P_{X,Y}$.*

- (d) [4 points] Change of Variable with Polar Coordinates. Let X and Y be random variables representing Cartesian coordinates in a plane. Let R and Θ be random variables representing the corresponding radius R and angle Θ in polar coordinates system. Then, Cartesian coordinates X and Y and polar coordinates are related by $X = R \cos(\Theta)$ and $Y = R \sin(\Theta)$. Using the change-of-variables formula, derive the joint probability density function $P_{X,Y}(r, \theta)$ in terms of $P_{R,\Theta}$.



There is a 1:1 mapping between R, Θ .

- (c) [2 points] Let $Z \sim \text{Uniform}[0, 3]$ and $X = Z^2$. Is this a valid normalization flow model?

No. Because given $X = 1$ there is no solution to Z . This is not an invertible map.

- (b) [2 points] Let $Z \sim \text{Uniform}[-2, 3]$ and $X = Z^2$. Is this a valid normalization flow model?

(a) [2 points] Let $Z \sim \text{Uniform}[-2, 3]$ (a uniform random variable over the interval $[-2, 3]$) and $X = Z^2$. What is $P_X(1)$?

$$P_X(1) = P_Z(f_\theta^{-1}(1)) \quad \text{where } f_\theta(z) = z^2$$

$$\frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$$

4. [10 points total] Normalizing Flow Models:

$Z \rightarrow x : f_\theta(z)$

8

$$\frac{\partial}{\partial \theta} \cos(\theta) = P_x(f_\theta(x)) \quad \text{for } \frac{\partial}{\partial \theta} \sin(\theta) = \cos(\theta)$$

5. [11 points total] **Variational Recurrent Neural Network**: Here we explore a recurrent version of the VAE for the purpose of modeling sequences. Drawing inspiration from simpler dynamic Bayesian networks (DBNs) such as HMMs and Kalman filters, this variational recurrent neural network (VRNN) augments a traditional RNN model with a sequence of latent variables; similar to a VAE, VRNN additionally models the dependencies between these latent random variables over time with another RNN.

A VRNN model defines a joint distribution $p(x_1, \dots, x_T, z_1, \dots, z_T)$ over two groups of random variables x_1, \dots, x_T and z_1, \dots, z_T . You can think of x_1, \dots, x_T as the elements of a time series, and z_1, \dots, z_T as latent variables that are also indexed by time. For brevity, we denote $\mathbf{x} \leq t = (x_1, \dots, x_t)$, $\mathbf{z} \leq t = (z_1, \dots, z_t)$, and we similarly define $\mathbf{z} \leq i$ and $\mathbf{z} \geq i$. In a VRNN, the joint is factored as

$$p_{\theta}(\mathbf{x} \leq T, \mathbf{z} \leq T) = \prod_{t=1}^T p_{\theta_1}(x_t | \mathbf{z} \leq t, \mathbf{x} \leq t) p_{\theta_2}(z_t | \mathbf{x} \leq t, \mathbf{z} \leq t)$$

where $p_{\theta_1}(x_i | \mathbf{z} \leq i, \mathbf{x} \leq i)$ and $p_{\theta_2}(z_i | \mathbf{x} \leq i, \mathbf{z} \leq i)$ are parameterized as RNNs, and $\theta = (\theta_1, \theta_2)$. As an edge case, $p_{\theta_1}(x_1 | \mathbf{z} \leq 1, \mathbf{x} \leq 1) = p_{\theta_1}(x_1 | z_1)$ and $p_{\theta_2}(z_1 | \mathbf{x} \leq 1, \mathbf{z} \leq 1) = p_{\theta_2}(z_1)$.

(a) [2 Points] How many parameters are needed to specify the model as a function of T ?

(c) [2 Points] Consider an alternative model

where $p_{\theta_1}(x_i | \mathbf{z} \leq i, \mathbf{x} \leq i)$ and $p_{\theta_2}(z_i | \mathbf{x} \leq i, \mathbf{z} \leq i)$ are parameterized as RNNs. Is this a valid autoregressive model?

No variable x_i needs to be conditioned on previous values to be used.

$$p_{\theta}(\mathbf{x} \leq T, \mathbf{z} \leq T) = \prod_{t=1}^T p_{\theta_1}(x_t | \mathbf{z} \leq t, \mathbf{x} \leq t) p_{\theta_2}(z_t | \mathbf{x} \leq t, \mathbf{z} \leq t)$$

(d) [4 Points] Learning: Because the marginal likelihood $p_{\theta}(\mathbf{x} \leq T)$ is intractable to compute, we need to resort to an ELBO for training. We consider the following encoder:

Show that the following ELBO for the VRNN described above

where $q_{\phi}(z_i | \mathbf{x} \leq i, \mathbf{z} \leq i)$ is parameterized using another RNN.

$$\log P_{\theta}(x_i | z_i) \geq E_{z \sim q_{\phi}(z_i | \mathbf{x} \leq i, \mathbf{z} \leq i)} \left[\log \frac{p_{\theta}(x_i | z_i, \mathbf{z} \leq i)}{q_{\phi}(z_i | \mathbf{x} \leq i, \mathbf{z} \leq i)} \right]$$

(d) [4 Points] Learning: Because the marginal likelihood $p_{\theta}(\mathbf{x} \leq T)$ is intractable to compute, we need to resort to an ELBO for training. We consider the following encoder:

Show that the following ELBO for the VRNN described above

where $q_{\phi}(z_i | \mathbf{x} \leq i, \mathbf{z} \leq i)$ is parameterized using another RNN.

$$\log P_{\theta}(x_i | z_i) \geq E_{z \sim q_{\phi}(z_i | \mathbf{x} \leq i, \mathbf{z} \leq i)} \left[\log \frac{p_{\theta}(x_i | z_i, \mathbf{z} \leq i)}{q_{\phi}(z_i | \mathbf{x} \leq i, \mathbf{z} \leq i)} \right]$$

(d) [4 Points] Learning: Because the marginal likelihood $p_{\theta}(\mathbf{x} \leq T)$ is intractable to compute, we need to resort to an ELBO for training. We consider the following encoder:

Show that the following ELBO for the VRNN described above

where $q_{\phi}(z_i | \mathbf{x} \leq i, \mathbf{z} \leq i)$ is parameterized using another RNN.

$$\log P_{\theta}(x_i | z_i) \geq E_{z \sim q_{\phi}(z_i | \mathbf{x} \leq i, \mathbf{z} \leq i)} \left[\log \frac{p_{\theta}(x_i | z_i, \mathbf{z} \leq i)}{q_{\phi}(z_i | \mathbf{x} \leq i, \mathbf{z} \leq i)} \right]$$

$$Z = \mu + \sigma e$$

~~Explain how the reparameterization trick~~

- (e) [1 Points] Assume the conditionals $q_\theta(z|x, \bar{z}, \bar{c})$ in the encoder RNN are Gaussians (with mean and variance computed by the RNN). Can the ELBO be optimized using the reparameterization trick?

where $D_{KL}(p||q)$ denotes KL divergence.

$$E_{\bar{z} \sim q_\theta(\bar{z}|x, \bar{x}, \bar{c})} \left[-D_{KL}(q_\theta(z|\bar{z}, \bar{x}, \bar{c}) || p_\theta(z|\bar{z})) \right]$$

is equivalent to

so the generators before were untrained

$$L = \frac{1}{T} \left[-D^{\phi}(G_{\theta}(z)) + D^{\phi}(G_{\theta}(z)) \right]$$

for non-saturating loss. The gradient of the generator is

This number is very good so $D^{\phi}(G_{\theta}(z)) \approx 0.9$

The non-saturating loss is more stable. In the beginning the

Hint: You might assume the generator produces samples that are easy to distinguish from real data at the beginning of training.

non-saturating loss results in more stable generator training and reduces the saturation problem.

(d) [4 Points] Advantage of non-saturating loss: Consider the derivatives of the saturating and

$$\frac{D^{\phi}(G_{\theta}(z))}{1 - D^{\phi}(G_{\theta}(z))}$$

(c) [2 points] Gradient of non-saturating Loss: What is the derivative of the non-saturating loss $L_{NS}(\theta)$ with respect to $D^{\phi}(G_{\theta}(z))$?

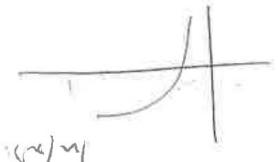
$$\frac{\partial L_{NS}(\theta)}{\partial D^{\phi}(G_{\theta}(z))} = \frac{1 - D^{\phi}(G_{\theta}(z))}{D^{\phi}(G_{\theta}(z))} \cdot \frac{\partial D^{\phi}(G_{\theta}(z))}{\partial \theta}$$

(b) [2 points] Gradient of the saturating Loss: What is the derivative of the saturating loss $L_S(\theta)$ with respect to $D^{\phi}(G_{\theta}(z))$?

then function requires $D^{\phi}(G_{\theta}(z))$ no large no small
that are equal to because finally to the min of each

$$\arg \min_{\theta} L_S(\theta) = \arg \min_{\theta} L_{NS}(\theta)$$

Show that for fixed z and D^{ϕ} :



$$L_{NS}(\theta) = -\log(D^{\phi}(G_{\theta}(z)))$$

Next, consider the following non-saturating loss:

$$L_S(\theta) = \log(1 - D^{\phi}(G_{\theta}(z)))$$

(a) [2 Points] Non-saturating Loss: Recall the original GAN generator loss for a single sample z :

Generative Adversarial Networks suffer from the saturation problem, which is when the generator G_{θ} cannot train as fast as the discriminator D^{ϕ} and stops learning. To overcome the saturation problem, a non-saturating loss is often used. Non-saturating loss is a subtle modification of the original generator loss, where the generator maximizes the probability of generated images being fake. than minimizing the probability of generated images being fake.

6. [10 Points total] GAN: Saturating and Non-Saturating Loss

(f) [4 points] KL Objective: Let's examine the KL more closely.

$\nabla \phi$

that we can evaluate tractably?

$$D_{KL}(p_T || q_\theta(x))$$

KL divergence

(e) [2 points] Learning from $p_T(x)$: Now suppose we want to train another autoregressive model $q_\theta(x)$ to learn the distribution of p_T . Is there an unbiased Monte Carlo estimate of the following

lower bound on the divergence of the two distributions?

(d) [2 points] Sampling from $p_T(x)$: Can we tractably sample from $p(x)$? Can we tractably sample

from $p_T(x)$ (for an arbitrary T)?

(c) [2 points] Temperature effect: Suppose we observe that samples x that have higher likelihood p_T generate higher quality samples (compared to p)

higher $p(x)$ are judged to have higher quality by human evaluators. How should we set T to have

$\nabla \phi$

to compute when n is large?

(b) [2 points] Partition function: Is the partition function $Z(2)$ for temperature $T = 2$ tractable

$$\log p_T(x) = \frac{1}{T} \log p(x) - \log Z(T)$$

$$\log Z(1) = \log p(x) - \log Z(1)$$

(a) [2 points] Partition function: What is the value of $Z(1)$, the partition function when $T = 1$?

where $p_T(x)$ is our temperature scaled joint distribution of interest and $Z(T)$ is the partition function.

$$\log p_T(x) = \frac{T}{1} \log p(x) - \log Z(T)$$

the following:

Let's consider a joint distribution $p(\mathbf{x})$ parameterized by an autoregressive model, where $\mathbf{x} = (x_1, x_2, \dots, x_n)$ is a collection of discrete random variables. We define joint temperature scaling by a scalar $T > 0$ as logits of individual variables is different from temperature scaling the logits of the joint distribution. calibrate the likelihood of the next token in language models. We showed that temperature scaling the probability of individual variables is different from temperature scaling the logits of the joint distribution.

In problem set 1, you were asked to implement temperature scaling, a commonly used technique to

7. [16 points total] Joint Temperature Scaling:

yes

(g) [2 points] Is there an unbiased estimator of Equation (*) that we can evaluate tractably?

$$\begin{aligned}
 & \left[\mathbb{E}_{\theta} \log q_{\theta} - \log \mathbb{E}_{\theta} q_{\theta} \right] = \mathbb{E}_{\theta} \left[\log \frac{q_{\theta}}{\mathbb{E}_{\theta} q_{\theta}} \right] \\
 & \left[\mathbb{E}_{\theta} \log q_{\theta} - \log \mathbb{E}_{\theta} q_{\theta} \right] = \int p_{\theta}(\mathbf{x}) \log \frac{p_{\theta}(\mathbf{x})}{\mathbb{E}_{\theta} p_{\theta}(\mathbf{x})} d\mathbf{x} \\
 & D_{KL}(p_{\theta} || \mathbb{E}_{\theta} p_{\theta})
 \end{aligned}$$

(Hint: Importance sampling.)

$$(*) \quad \mathbb{E}_{\theta} \left[\log \frac{q_{\theta}(\mathbf{x})}{\mathbb{E}_{\theta} q_{\theta}(\mathbf{x})} \right]$$

minimizes the following:

Show that finding the θ that minimizes $D_{KL}(p_{\theta} || q_{\theta}(\mathbf{x}))$ is equivalent to finding the θ that

