

My Notes for XCS236

Fangfei Li

December 15, 2025

Contents

1	Introduction & Background	1
1.1	Introduction	2
1.2	Background	6
2	Autoregressive Models	9
2.1	Introduction	10
2.2	Examples of AR Models	11
2.3	Fully Visible Sigmoid Belief Network (FVSBN)	14
2.4	Neural Autoregressive Density Estimation (NADE)	15
2.5	General Discrete Distributions	16
2.6	Real-valued Neural Autoregressive Density Estimator (RNADE)	17
2.7	Other	18
3	Maximum Likelihood Learning (MLE)	23
3.1	Introduction	24
3.2	MLE for Autoregressive Models	28
3.3	Bias and Variance	30
4	Latent Variable Models (VAEs)	31
4.1	Latent Variable Models	32
4.2	Notations of VAEs	34
4.3	Evidence Lower Bound (ELBO)	35
4.4	β -VAE	38
4.5	Importance Weighted Autoencoder (IWAE) Bound	39
4.6	Assumptions of VAEs	42
4.7	Workflow of VAEs	43
5	Normalizing Flows	45
5.1	Introduction	46
5.2	Comparison of VAE & NF	47
5.3	Normalizing Flows	48
6	Generative Adversarial Networks (GANs)	50
6.1	Motivation of Likelihood-Free Learning	50
6.2	Notation for GAN	51

6.3	GAN	52
6.4	Discriminator Loss	53
6.5	Generator Minimax Loss	55
6.6	Generator Non-Saturating Loss	56
6.7	Generator KL-based Loss	57
6.8	WGAN Critic Loss	58
6.9	WGAN-GP Critic Loss	58
6.10	WGAN Generator Loss	58
6.11	Standard GAN	59
6.12	Non-saturating GAN	61
6.13	KL-based GAN	62
6.14	Conditional GAN with Projection Discriminator	63
6.15	Wasserstein GAN (WGAN)	64
6.16	Wasserstein GAN with Gradient Penalty (WGAN-GP)	65
6.17	Issue Analysis	66
6.18	f -divergence	68

1 Introduction & Background

1.1 Introduction

Definition: Generative model (or statistical generative model) is,

- Probability distribution: $p(\mathbf{x})$ with $\mathbf{x} \in \mathcal{X}$
- Or in some cases, conditional probability distribution: $p(\mathbf{x} | \mathbf{c})$ with $\mathbf{x} \in \mathcal{X}$ and $\mathbf{c} \in \mathcal{C}$

Type: Generative model includes:

- Probabilistic graphical models (PGM):
 - Directed Probabilistic graphical models (DPGM): Bayesian Networks (BN)
 - Undirected Probabilistic graphical models (UPGM): Markov random field (MRF)
- Deep generative models (DGM):
 - Density-based generative models: Explicit $p(\mathbf{x})$, trained by maximum likelihood.
 - * Autoregressive models (ARM)
 - * Variational autoencoders (VAE)
 - * Normalizing flows
 - Likelihood-free generative models: Implicit $p(\mathbf{x})$, generate samples from noise. through a generator network
 - * Energy based models (EBM): learn from energy differences
 - * Diffusion models: learn from denoising objectives
 - * Generative adversarial networks (GAN): learn from sample comparisons via a discriminator

Applications: Generative model can be used for,

1. Generating data:
 - Generate new data samples (e.g., image, text, audio, etc).
 - Create synthetic data for training discriminative models.
 - Usually use a conditional generative model $p(\mathbf{x} | \mathbf{c})$ to control what is generated.
 - (a) text-to-image model: \mathbf{c} = text prompt, \mathbf{x} = image
 - (b) Image-to-text model: \mathbf{c} = image, \mathbf{x} = text (image captioning)
 - (c) image-to-image model: \mathbf{c} = image, \mathbf{x} = image (image colorization, inpainting, uncropping, JPEG artefact restoration)
 - (d) speech-to-text model: \mathbf{c} = sequence of sounds, \mathbf{x} = sequence of words (automatic speech recognition (ASR))
 - (e) sequence-to-sequence
 - \mathbf{c} = sequence of English words, \mathbf{x} = sequence of French words (machine translation)
 - \mathbf{c} = initial prompt, \mathbf{x} = continuation of the text (text generation)
 - Difference between discriminative model and conditional generative model
 - (a) Discriminative (predictive) model: A function $f_\theta : \mathbf{x} \rightarrow \mathbf{y}$, where there is a single correct output \mathbf{y} for each input \mathbf{x} .

(b) Conditional generative model: A probability distribution $p(\mathbf{x} \mid \mathbf{c})$, where multiple outputs \mathbf{x} may correspond to the same input \mathbf{c} . This makes generative models more harder to evaluate.

2. Density estimation:

- Evaluate the probability of an observed data vector, i.e., compute $p(\mathbf{x})$.
- Used in outlier detection, data compression, generative classifiers, model comparison, etc.
- Approaches for outlier detection:
 - (a) kernel density estimation (KDE): a simple non-parametric generative model only for low-dimensional data. For higher-dimensional data, it suffers from the curse of dimensionality, making density estimation inaccurate and computationally expensive. Need parametric generative models (e.g., Gaussian mixtures, Hidden Markov Models).
 - (b) deep generative model: can learn complex, high-dimensional data distributions.

3. Imputation:

- “Filling in” missing values of a data vector or data matrix
- Used in in-painting occluded pixels in an image
- Approaches for imputation
 - (a) mean value imputation
 - (b) multiple imputation

4. Structure discovery:

- For generative models with latent variables, the latent vector \mathbf{z} is assumed to represent the underlying “causes” that generate the observed data \mathbf{x} . By Bayes’ rule, the posterior distribution is obtained as

$$p(\mathbf{z} \mid \mathbf{x}) = \frac{p(\mathbf{z})p(\mathbf{x} \mid \mathbf{z})}{p(\mathbf{x})} \propto p(\mathbf{z})p(\mathbf{x} \mid \mathbf{z})$$

which is useful for discovering latent, low-dimensional patterns in the data.

5. Latent space interpolation:

- Certain latent variable models can generate new samples with desired semantic properties by interpolating between existing datapoints in the latent space.
 - (a) Given datapoints \mathbf{x}_1 and \mathbf{x}_2 , obtain their latent representation using the encoder $e(\cdot)$, $\mathbf{z}_1 = e(\mathbf{x}_1)$ and $\mathbf{z}_2 = e(\mathbf{x}_2)$.
 - (b) Construct an interpolated latent point, $\mathbf{z} = \lambda\mathbf{z}_1 + (1 - \lambda)\mathbf{z}_2$.
 - (c) Decoding \mathbf{z} using the decoder $d(\cdot)$, $\mathbf{x} = d(\mathbf{z})$. \mathbf{x} typically combines semantic features from both \mathbf{x}_1 and \mathbf{x}_2 .
 - (d) Linear interpolation is commonly used because the learned latent manifold often exhibits approximately zero curvature in practice. However, in some cases nonlinear interpolation (e.g., spherical interpolation, geodesic interpolation) may better respect the geometry of the latent space.

6. Latent space arithmetic:

- Increase or decrease the amount of a desired “semantic factor of variation”.

7. Generative design:

- Generate candidate objects (e.g., molecules) which have desired properties.
- Fit a VAE to unlabeled samples, then to perform Bayesian optimization.

8. Model-based reinforcement learning:

9. Representation learning:

- Learn latent features z from the observed data x , that can be used for downstream supervised tasks.

10. Data compression:

Slide Note

- Make sense of complex, high-dimensional data (e.g., Image, speech signal, a sequence of tokens/characters, etc.) by converting them into useful representations that reveal their true meaning.
- Quote:
 - Richard Feynman: “What I cannot create, I do not understand.”
 - Philosophy Behind Generative modeling: “What I understand, I can create.”
- Model structure:
 - Generation (graphic): High level description → Raw sensory output
 - Inference (inverse graphic): Raw sensory output → High level description

Note: Many of models in this course will have similar structure: Generation + Inference
- A statistical generative model is a probability distribution $p(x)$ that is learned from,
 - Data (samples) - more weight
 - Prior Knowledge - less weight
- Data → Data simulator (= Statistical model = Generative model = Probability distribution $p(x)$)
 - Generation (Sampling): Called generative because sampling from $p(x)$ generates new sample.
Control signal (e.g., texts/captions, images, medical signals, etc.) → Data simulator → New data
 - Density estimation: Used for tasks such as outlier or anomaly detection.
New data → Data simulator → Density

Model Type: Models can be categorized along two dimensions:

- What the model learns
 1. Generative models
 - Learn the underlying data distribution and can generate new samples.
 - E.g., Generative Pretrained Transformer (GPT)
 2. Discriminative models
 - Learn to predict labels or distinguish classes, modeling only $p(y | x)$.
 - E.g., Bidirectional Encoder Representations from Transformers (BERT)
- How the model is trained
 1. Supervised learning
 - Trained with labeled input–output pairs.
 - E.g., classifiers
 2. Self-supervised learning
 - Trained using labels derived automatically from the data itself.
 - E.g., GPT (predicts next token), BERT (predicts masked tokens)
 3. Unsupervised learning

- Trained on unlabeled data to discover structure or generate data.
- E.g., Variational Autoencoder (VAE)

1.2 Background

Examples of Data Distributions

- Bernoulli distribution

– $X \sim \text{Bern}(p)$:

$$P(X = 1) = p, \quad P(X = 0) = 1 - p, \quad \text{Val}(X) = \{0, 1\}$$

– A black–white image with n pixels:

1. Each pixel X_i is Bernoulli with $\text{Val}(X_i) = \{0, 1\}$.
2. The joint distribution of the full image is $p(x_1, x_2, \dots, x_n)$.
3. The number of possible images (states) is 2^n .
4. The number of parameters required to specify the joint distribution is $2^n - 1$.

- Categorical distribution

– $Y \sim \text{Categorical}(p_1, p_2, \dots, p_m)$:

$$P(Y = i) = p_i, \quad \sum_{i=1}^m p_i = 1, \quad \text{Val}(Y) = \{1, 2, \dots, m\}$$

– Color of a single RGB pixel:

1. Each channel (R, G, B) is categorical with $\{0, \dots, 255\}$.
2. The full pixel has a joint distribution $p(r, g, b)$.
3. The total number of possible pixel values is 256^3 .
4. The number of parameters required to specify the joint distribution is $256^3 - 1$.

Independence vs. Conditional Independence

1. Independence assumption is too strong.
2. Conditional independence assumption is better.

Chain Rule

$$p(x_1, \dots, x_n) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2) \cdots p(x_n | x_1, \dots, x_{n-1})$$

Bayes' Rule

$$p(x | y) = \frac{p(y | x) p(x)}{p(y)}$$

$$p(y) = \sum_x p(y | x) p(x)$$

Bayesian Network (BN):

- A BN is defined by a directed acyclic graph (DAG).
- The DAG structure encodes conditional independence relationships among variables.

- The joint distribution factorizes into local conditional probability distributions (CPDs),

$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | \mathbf{x}_{\text{Parent}(i)})$$

- This factorization defines a valid probability distribution (i.e., it sums to 1) because the DAG provides an ordering consistent with the conditional dependencies.

Neural Models vs. Bayesian Networks:

- Logistic regression is generally more powerful than naive Bayes for predicting Y because it does not rely on the conditional independence assumption $X_i \perp \mathbf{X}_{-i} | Y$. A sufficiently deep neural network has greater flexibility, as it can approximate arbitrary functions.
- Naive Bayes, however, is more versatile than logistic regression in certain settings. It can be used for tasks beyond predicting Y . Even when the goal is prediction, naive Bayes can still operate when some features in \mathbf{X} are missing, while logistic regression typically requires all features to be observed.

Generative vs. Discriminative Models

- A joint distribution can be factorized in two equivalent ways:

- Generative form: $p(y, \mathbf{x}) = p(\mathbf{x} | y)p(y)$
- Discriminative form: $p(y, \mathbf{x}) = p(y | \mathbf{x})p(\mathbf{x})$

- Examples of generative models:

- Naive Bayes (Naive BN): Under the conditional independence assumption $X_i \perp \mathbf{X}_{-i} | Y$, the joint distribution factorizes as $p(y, x_1, \dots, x_n) = p(y) \prod_{i=1}^n p(x_i | y)$. The posterior becomes

$$p(Y = 1 | x_1, \dots, x_n) = \frac{p(Y = 1) \prod_{i=1}^n p(x_i | Y = 1)}{p(Y = 0) \prod_{i=1}^n p(x_i | Y = 0) + p(Y = 1) \prod_{i=1}^n p(x_i | Y = 1)}$$

- Example of discriminative models

- Logistic regression (linear dependence): Let sigmoid function $\sigma(z) = \frac{1}{1+e^{-z}}$ (used for generating probability with values in $[0, 1]$). Given parameters $\boldsymbol{\alpha} = (\alpha_0, \alpha_1, \dots, \alpha_n)$,

$$p_{\text{logit}}(Y = 1 | \mathbf{x}; \boldsymbol{\alpha}) = \sigma(\alpha_0 + \sum_{i=1}^n \alpha_i x_i)$$

- Neural network model (non-linear dependence): Let $\mathbf{h} = f(A\mathbf{x} + \mathbf{b})$ be a nonlinear hidden representation. Then

$$p_{\text{Neural}}(Y = 1 | \mathbf{x}; \boldsymbol{\alpha}, A, \mathbf{b}) = \sigma(\alpha_0 + \sum_{i=1}^n \alpha_i h_i)$$

Continuous Variable

Continuous variables

- If X is a continuous random variable, we can usually represent it using its **probability density function** $p_X : \mathbb{R} \rightarrow \mathbb{R}^+$. However, we cannot represent this function as a table anymore. Typically consider parameterized densities:
 - Gaussian: $X \sim \mathcal{N}(\mu, \sigma)$ if $p_X(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$
 - Uniform: $X \sim \mathcal{U}(a, b)$ if $p_X(x) = \frac{1}{b-a} \mathbf{1}[a \leq x \leq b]$
 - Etc.
- If \mathbf{X} is a continuous random vector, we can usually represent it using its **joint probability density function**:
 - Gaussian: if $p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}))$
 - Chain rule, Bayes rule, etc all still apply. For example,

$$p_{X,Y,Z}(x, y, z) = p_X(x)p_{Y|X}(y | x)p_{Z|{X,Y}}(z | x, y)$$

• CS236, Stanford University Deep Generative Models Lecture 2 28 / 29

Continuous variables

- This means we can still use Bayesian networks with continuous (and discrete) variables. Examples:
- Mixture of 2 Gaussians:** Bayes net $Z \rightarrow X$ with factorization $p_{Z,X}(z, x) = p_Z(z)p_{X|Z}(x | z)$ and
 - $Z \sim \text{Bernoulli}(p)$
 - $X | (Z=0) \sim \mathcal{N}(\mu_0, \sigma_0)$, $X | (Z=1) \sim \mathcal{N}(\mu_1, \sigma_1)$
 - The parameters are $p, \mu_0, \sigma_0, \mu_1, \sigma_1$
- Bayes net $Z \rightarrow X$ with factorization $p_{Z,X}(z, x) = p_Z(z)p_{X|Z}(x | z)$
 - $Z \sim \mathcal{U}(a, b)$
 - $X | (Z=z) \sim \mathcal{N}(z, \sigma)$
 - The parameters are a, b, σ
- Variational autoencoder:** Bayes net $Z \rightarrow X$ with factorization $p_{Z,X}(z, x) = p_Z(z)p_{X|Z}(x | z)$ and
 - $Z \sim \mathcal{N}(0, 1)$
 - $X | (Z=z) \sim \mathcal{N}(\mu_\theta(z), e^{\sigma_\phi(z)})$ where $\mu_\theta : \mathbb{R} \rightarrow \mathbb{R}$ and σ_ϕ are neural networks with parameters (weights) θ, ϕ respectively
 - Note:** Even if μ_θ, σ_ϕ are very deep (flexible), functional form is still Gaussian

• CS236, Stanford University Deep Generative Models Lecture 2 29 / 29

2 Autoregressive Models

2.1 Introduction

Two-Step Process

1. Model representation: Parameterize a model family $\{p_\theta(x), \theta \in \Theta\}$ to represent $p(x)$
2. Model Learning: Estimate (or optimize) the parameters θ using the training data \mathcal{D}

Definition: Autoregressive (AR) models use parameterized functions to predict each data element (e.g., a pixel) based on all previously observed elements. The joint distribution is defined as

$$p(x_1, \dots, x_d) = \prod_{i=1}^d p(x_i | x_1, \dots, x_{i-1}) = \prod_{i=1}^d p(x_i | \mathbf{x}_{<i})$$

It's self-supervised learning, since the model learns from the data itself without manual labels.

2.2 Examples of AR Models

Example 1 – Recurrent Neural Networks (RNNs): RNNs are autoregressive (AR) models that process a sequence of input vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$, with each $\mathbf{x}_t \in \mathbb{R}^I$. At each time step, the network maintains a hidden state $\mathbf{h}_t \in \mathbb{R}^H$.

- Dimensions:
 - I : input dimension
 - H : hidden state dimension (hyperparameter)
 - V : output dimension
- Parameters:
 - $W_x \in \mathbb{R}^{H \times I}$: input-to-hidden weight matrix
 - $W_h \in \mathbb{R}^{H \times H}$: recurrent (hidden-to-hidden) weight matrix
 - $\mathbf{b}_h \in \mathbb{R}^H$: bias for the hidden state
 - $W_o \in \mathbb{R}^{V \times H}$: hidden-to-output weight matrix
 - $\mathbf{b}_o \in \mathbb{R}^V$: output bias
- Hidden state update using a nonlinear activation $f(\cdot)$ (e.g., tanh, ReLU)

$$\mathbf{h}_t = f(W_h \mathbf{h}_{t-1} + W_x \mathbf{x}_t + \mathbf{b}_h)$$

- Output distribution:

$$p(\mathbf{x}_{t+1} | \mathbf{x}_{\leq t}) = g(W_o \mathbf{h}_t + \mathbf{b}_o)$$

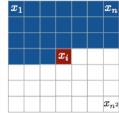
where $g(\cdot)$ is typically a softmax.

- $W_o \mathbf{h}_t + \mathbf{b}_o$: logits, which are the raw, unnormalized scores before applying softmax to convert them into probabilities.

Example 2 – Long Short-Term Memory Networks (LSTMs): LSTMs are AR models, which is an improved form of RNNs that capture long-range dependencies through gating mechanisms.

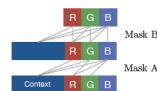
Example 3 – PixelRNN: PixelRNNs are AR models for images that generate each pixel conditioned on all previously generated pixels.

Pixel RNN (Oord et al., 2016)



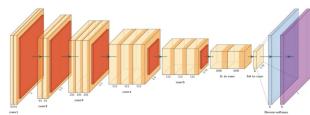
- ① Model images pixel by pixel using raster scan order
- ② Each pixel conditional $p(x_t | x_{1:t-1})$ needs to specify 3 colors

$$p(x_t | x_{1:t-1}) = p(x_t^{red} | x_{1:t-1})p(x_t^{green} | x_{1:t-1}, x_t^{red})p(x_t^{blue} | x_{1:t-1}, x_t^{red}, x_t^{green})$$
and each conditional is a categorical random variable with 256 possible values
- ③ Conditionals modeled using RNN variants. LSTMs + masking (like MADE)



• Stefano Ermon (AI Lab) Deep Generative Models Lecture 3 28 / 36

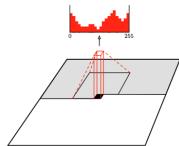
Convolutional Architectures



Convolutions are natural for image data and easy to parallelize on modern hardware.

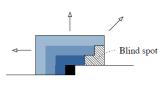
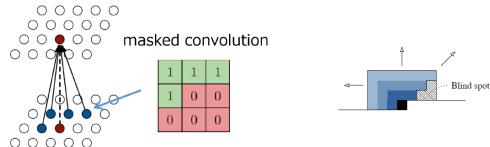
• Stefano Ermon (AI Lab) Deep Generative Models Lecture 3 30 / 36

PixelCNN (Oord et al., 2016)



Idea: Use convolutional architecture to predict next pixel given context (a neighborhood of pixels).

Challenge: Has to be autoregressive. Masked convolutions preserve raster scan order. Additional masking for colors order.



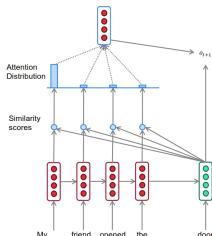
• Stefano Ermon (AI Lab) Deep Generative Models Lecture 3 31 / 36

Example 4 – Generative Pre-trained Transformer (GPT): GPT models are transformer-based autoregressive

models that predict the next token given all preceding tokens.

Attention Based Models: A broad class of neural models that use attention mechanisms to selectively weight input information.

Attention based models



Attention mechanism to compare a *query* vector to a set of *key* vectors

- ❶ Compare current hidden state (*query*) to all past hidden states (*keys*), e.g., by taking a dot product
- ❷ Construct attention distribution to figure out what parts of the history are relevant, e.g., via a softmax
- ❸ Construct a summary of the history, e.g., by weighted sum
- ❹ Use summary and current hidden state to predict next token/word

Generative Transformers:

- Generative transformers = Transformers used for generative modeling
- Examples: GPT

Generative Transformers



Current state of the art (GPTs): replace RNN with Transformer

- Attention mechanisms to adaptively focus only on relevant context
- Avoid recursive computation. Use only self-attention to enable parallelization
- Needs **masked** self-attention to preserve autoregressive structure
- Demo: <https://transformer.huggingface.co/doc/gpt2-large>
- Demo: <https://huggingface.co/spaces/huggingface-projects/llama-2-13b-chat>

2.3 Fully Visible Sigmoid Belief Network (FVSBN)

Definition: In particular, an AR model is called a Fully Visible Sigmoid Belief Network when

1. Each variable is binary (Bernoulli).
2. Each conditional probability is modeled as a sigmoid function of a weighted sum of all previous variables.

A FVSBN is a probabilistic graphical model where all variables are observed (fully visible) and connected in a DAG according to an autoregressive ordering.

Example: For a black-white image with 28×28 pixels (so $d = 784$), the joint distribution can be written autoregressively as

$$p(x_1, \dots, x_{784}) = p_{\text{CPT}}(x_1; \alpha^1)p_{\text{logit}}(x_2 | x_1, \alpha^2) \cdots p_{\text{logit}}(x_d | x_1, \dots, x_{d-1}, \alpha^d),$$

Here, each conditional variables $X_i | x_1, \dots, x_{i-1}$ is Bernoulli with parameter

$$\begin{aligned} \hat{x}_i &= p(X_i = 1 | x_1, \dots, x_{i-1}; \alpha^i) \\ &= p(X_i = 1 | x_{<i}; \alpha^i) \\ &= \sigma(\alpha_0^i + \sum_{j=1}^{i-1} \alpha_j^i x_j) \end{aligned}$$

and $\alpha^i = (\alpha_0^i, \alpha_1^i, \dots, \alpha_{i-1}^i)$.

Thus the conditional distributions are simply

$$\begin{aligned} p_{\text{CPT}}(X_1 = 1; \alpha^1) &= \alpha^1 \\ p_{\text{CPT}}(X_1 = 0; \alpha^1) &= 1 - \alpha^1 \\ p_{\text{logit}}(X_2 = 1 | x_1; \alpha^2) &= \sigma(\alpha_0^2 + \alpha_1^2 x_1) \\ p_{\text{logit}}(X_2 = 0 | x_1; \alpha^2) &= 1 - \sigma(\alpha_0^2 + \alpha_1^2 x_1) \\ &\dots \\ p_{\text{logit}}(X_d = 1 | x_1, \dots, x_{d-1}; \alpha^d) &= \sigma(\alpha_0^d + \sum_{j=1}^{d-1} \alpha_j^d x_j) \\ p_{\text{logit}}(X_d = 0 | x_1, \dots, x_{d-1}; \alpha^d) &= 1 - \sigma(\alpha_0^d + \sum_{j=1}^{d-1} \alpha_j^d x_j) \end{aligned}$$

The total number of parameters is $1 + 2 + \dots + d \approx \frac{d^2}{2}$.

2.4 Neural Autoregressive Density Estimation (NADE)

Neural Autoregressive Density Estimation extends the FVSBN model by replacing the linear component of the logistic regression with a single-layer neural network, while still maintaining a sigmoid output layer. Specifically,

$$\begin{aligned}\hat{x}_i &= p(X_i = 1 \mid x_1, \dots, x_{i-1}; A_i, \mathbf{c}_i, \boldsymbol{\alpha}_i, b_i) \\ &= p(X_i = 1 \mid \mathbf{x}_{<i}; A_i, \mathbf{c}_i, \boldsymbol{\alpha}_i, b_i) \\ &= \sigma(\boldsymbol{\alpha}_i \mathbf{h}_i + b_i)\end{aligned}$$

Here,

1. $\mathbf{x}_{<i} \in \mathbb{R}^{(i-1) \times 1}$ is vector of previous variables x_1, \dots, x_{i-1} .
2. $A_i \in \mathbb{R}^{H \times (i-1)}$ is weight matrix connecting inputs $\mathbf{x}_{<i}$ to hidden units.
3. $\mathbf{c}_i \in \mathbb{R}^{H \times 1}$ is bias vector for hidden layer.
4. $\mathbf{h}_i \in \mathbb{R}^{H \times 1}$ is hidden layer activations, i.e.,

$$\mathbf{h}_i = \sigma(A_i \mathbf{x}_{<i} + \mathbf{c}_i) = \sigma(A_i \begin{pmatrix} x_1 \\ \vdots \\ x_{i-1} \end{pmatrix} + \mathbf{c}_i)$$

5. $\boldsymbol{\alpha}_i \in \mathbb{R}^{1 \times H}$ is output weight vector.
6. $b_i \in \mathbb{R}$ is output bias.

Instead of having a separate weight matrix A_i for each i , one can use a shared weight matrix $W \in \mathbb{R}^{H \times d}$. For each i , only the first $i - 1$ columns are used. And also shared $\mathbf{c} \in \mathbb{R}^{H \times 1}$ is bias vector for hidden layer.

$$\mathbf{h}_i = \sigma(W_{:, <i} \mathbf{x}_{<i} + \mathbf{c})$$

NADE: Neural Autoregressive Density Estimation

- Tie weights to *reduce the number of parameters* and *speed up computation* (see blue dots in the figure):

$$\mathbf{h}_i = \sigma(W_{:, <i} \mathbf{x}_{<i} + \mathbf{c})$$

$$\hat{x}_i = p(x_i \mid x_1, \dots, x_{i-1}) = \sigma(\boldsymbol{\alpha}_i \mathbf{h}_i + b_i)$$
- For example

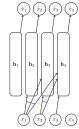
$$\mathbf{h}_2 = \sigma \left(\underbrace{\begin{pmatrix} \vdots \\ \mathbf{w}_1 \\ \vdots \end{pmatrix}}_{W_{:, <2}} \mathbf{x}_1 + \mathbf{c} \right)$$

$$\mathbf{h}_3 = \sigma \left(\underbrace{\begin{pmatrix} \vdots \\ \mathbf{w}_1 \mathbf{w}_2 \\ \vdots \end{pmatrix}}_{W_{:, <3}} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} \vdots \\ \mathbf{c} \end{pmatrix} \right)$$

$$\mathbf{h}_4 = \sigma \left(\underbrace{\begin{pmatrix} \vdots \\ \mathbf{w}_1 \mathbf{w}_2 \mathbf{w}_3 \\ \vdots \end{pmatrix}}_{W_{:, <4}} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} + \begin{pmatrix} \vdots \\ \mathbf{c} \end{pmatrix} \right)$$
- If $\mathbf{h}_i \in \mathbb{R}^d$, how many total parameters? Linear in n : weights $W \in \mathbb{R}^{H \times d}$, biases $\mathbf{c} \in \mathbb{R}^d$, and n logistic regression coefficient vectors $\boldsymbol{\alpha}_i, b_i \in \mathbb{R}^{d+1}$. Probability is evaluated in $O(nd)$.

2.5 General Discrete Distributions

General discrete distributions



How to model non-binary discrete random variables $X_i \in \{1, \dots, K\}$? E.g., pixel intensities varying from 0 to 255

One solution: Let \hat{x}_i parameterize a categorical distribution

$$\begin{aligned}\mathbf{h}_i &= \sigma(W_{:,i} \mathbf{x}_{<i} + \mathbf{c}) \\ p(x_i | x_1, \dots, x_{i-1}) &= \text{Cat}(p_i^1, \dots, p_i^K) \\ \hat{\mathbf{x}}_i &= (p_i^1, \dots, p_i^K) = \text{softmax}(A_i \mathbf{h}_i + \mathbf{b}_i)\end{aligned}$$

Softmax generalizes the sigmoid/logistic function $\sigma(\cdot)$ and transforms a vector of K numbers into a vector of K probabilities (non-negative, sum to 1).

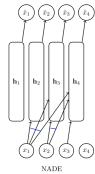
$$\text{softmax}(\mathbf{a}) = \text{softmax}(a^1, \dots, a^K) = \left(\frac{\exp(a^1)}{\sum_i \exp(a^i)}, \dots, \frac{\exp(a^K)}{\sum_i \exp(a^i)} \right)$$

In numpy: `np.exp(a)/np.sum(np.exp(a))`



2.6 Real-valued Neural Autoregressive Density Estimator (RNADE)

RNADE



How to model continuous random variables $X_i \in \mathbb{R}$? E.g., speech signals

Solution: let \hat{x}_i parameterize a continuous distribution

E.g., In a mixture of K Gaussians,

$$p(x_i | x_1, \dots, x_{i-1}) = \sum_{j=1}^K \frac{1}{K} \mathcal{N}(x_i; \mu_i^j, \sigma_i^j)$$

$$\mathbf{h}_i = \sigma(W_{:, < i} \mathbf{x}_{< i} + \mathbf{c})$$

$$\hat{\mathbf{x}}_i = (\mu_i^1, \dots, \mu_i^K, \sigma_i^1, \dots, \sigma_i^K) = f(\mathbf{h}_i)$$

$\hat{\mathbf{x}}_i$ defines the mean and standard deviation of each of the K Gaussians (μ_i^j, σ_i^j) .

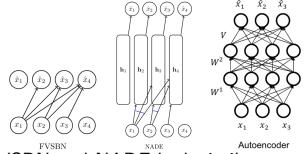
Can use exponential $\exp(\cdot)$ to ensure non-negativity



2.7 Other

Autoregressive Models vs. Autoencoders

Autoregressive models vs. autoencoders



- On the surface, FVSBN and NADE look similar to an **autoencoder**:
- an **encoder** $e(\cdot)$. E.g., $e(x) = \sigma(W^2(W^1x + b^1) + b^2)$
- a **decoder** such that $d(e(x)) \approx x$. E.g., $d(h) = \sigma(Vh + c)$.
- Loss function for dataset \mathcal{D}

$$\begin{aligned} \text{Binary r.v.: } & \min_{W^1, W^2, b^1, b^2, V, c} \sum_{x \in \mathcal{D}} \sum_i -x_i \log \hat{x}_i - (1 - x_i) \log(1 - \hat{x}_i) \\ \text{Continuous r.v.: } & \min_{W^1, W^2, b^1, b^2, V, c} \sum_{x \in \mathcal{D}} \sum_i (x_i - \hat{x}_i)^2 \end{aligned}$$

- e and d are constrained so that we don't learn identity mappings. Hope that $e(x)$ is a meaningful, compressed representation of x (feature learning)
- A vanilla autoencoder is *not* a generative model: it does not define a distribution over x we can sample from to generate new data points.



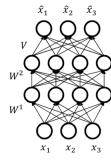
Stefano Ermon (AI Lab)

Deep Generative Models

Lecture 3

15 / 36

Autoregressive autoencoders



- On the surface, FVSBN and NADE look similar to an **autoencoder**. Can we get a generative model from an autoencoder?
- We need to make sure it corresponds to a valid Bayesian Network (DAG structure), i.e., we need an *ordering* for chain rule. If ordering is 1, 2, 3, then:
 - \hat{x}_1 cannot depend on any input $x = (x_1, x_2, x_3)$. Then at generation time we don't need any input to get started
 - \hat{x}_2 can only depend on x_1
 - ...
- Bonus:** we can use a single neural network (with n inputs and outputs) to produce all the parameters \hat{x} in a single pass. In contrast, NADE requires n passes. Much more efficient on modern hardware.



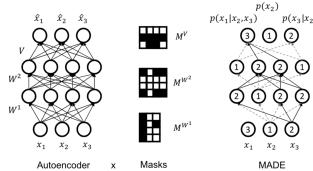
Stefano Ermon (AI Lab)

Deep Generative Models

Lecture 3

16 / 36

MADE: Masked Autoencoder for Distribution Estimation

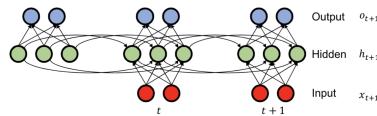


- ❶ **Challenge:** An autoencoder that is autoregressive (DAG structure)
- ❷ **Solution:** use masks to disallow certain paths (Germain et al., 2015). Suppose ordering is x_2, x_3, x_1 , so $p(x_1, x_2, x_3) = p(x_2)p(x_3 | x_2)p(x_1 | x_2, x_3)$.
 - ❸ The unit producing the parameters for $\hat{x}_2 = p(x_2)$ is not allowed to depend on any input. Unit for $p(x_3 | x_2)$ only on x_2 . And so on...
 - ❹ For each unit in a hidden layer, pick a random integer i in $[1, n - 1]$. That unit is allowed to depend only on the first i inputs (according to the chosen ordering).
 - ❺ Add mask to preserve this invariant: connect to all units in previous layer with smaller or equal assigned number (strictly $<$ in final layer)

Autoregressive Models vs. RNNs

RNN: Recurrent Neural Nets

Challenge: model $p(x_t | x_{1:t-1}; \alpha^t)$. "History" $x_{1:t-1}$ keeps getting longer.
Idea: keep a summary and recursively update it



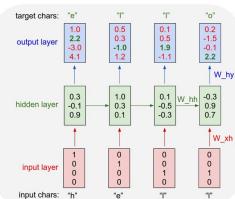
$$\text{Summary update rule: } h_{t+1} = \tanh(W_{hh}h_t + W_{xh}x_{t+1})$$

$$\text{Prediction: } o_{t+1} = W_{hy}h_{t+1}$$

$$\text{Summary initialization: } h_0 = b_0$$

- ❶ Hidden layer h_t is a summary of the inputs seen till time t
- ❷ Output layer o_{t-1} specifies parameters for conditional $p(x_t | x_{1:t-1})$
- ❸ Parameterized by b_0 (initialization), and matrices W_{hh}, W_{xh}, W_{hy} . Constant number of parameters w.r.t n !

Example: Character RNN (from Andrej Karpathy)



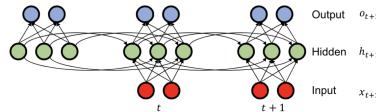
- ❶ Suppose $x_i \in \{h, e, l, o\}$. Use one-hot encoding:
 - ❷ h encoded as $[1, 0, 0, 0]$, e encoded as $[0, 1, 0, 0]$, etc.
- ❸ **Autoregressive:** $p(x = hello) = p(x_1 = h)p(x_2 = e | x_1 = h)p(x_3 = l | x_1 = h, x_2 = e) \cdots p(x_5 = o | x_1 = h, x_2 = e, x_3 = l, x_4 = l)$
- ❹ For example,

$$p(x_2 = e | x_1 = h) = \text{softmax}(o_1) = \frac{\exp(2.2)}{\exp(1.0) + \cdots + \exp(4.1)}$$

$$o_1 = W_{hy}h_1$$

$$h_1 = \tanh(W_{hh}h_0 + W_{xh}x_1)$$

RNN: Recurrent Neural Nets



Pros:

- ➊ Can be applied to sequences of arbitrary length.
- ➋ Very general: For every computable function, there exists a finite RNN that can compute it

Cons:

- ➌ Still requires an ordering
- ➍ Sequential likelihood evaluation (very slow for training)
- ➎ Sequential generation (unavoidable in an autoregressive model)

• Stefano Ermon (AI Lab) Deep Generative Models Lecture 3 20 / 36

Example: Character RNN (from Andrej Karpathy)

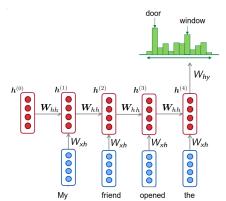
Train 3-layer RNN with 512 hidden nodes on all the works of Shakespeare.
Then sample from the model:

KING LEAR: O, if you were a feeble sight, the courtesy of your law,
Your sight and several breath, will wear the gods
With his heads, and my hands are wonder'd at the deeds,
So drop upon your lordship's head, and your opinion
Shall be against your honour.

Note: generation happens **character by character**. Needs to learn valid words, grammar, punctuation, etc.

• Stefano Ermon (AI Lab) Deep Generative Models Lecture 3 21 / 36

Issues with RNNs



Issues with RNN models

- ➊ A single hidden vector needs to summarize all the (growing) history.
For example, $h^{(4)}$ needs to summarize the meaning of "My friend opened the".
- ➋ Sequential evaluation, cannot be parallelized
- ➌ Exploding/vanishing gradients when accessing information from many steps back

• Stefano Ermon (AI Lab) Deep Generative Models Lecture 3 25 / 36

RNNs can be autoregressive. E.g., when generating text, RNNs use the previous token to predict the next. That

makes them autoregressive in behavior. Transformers (like GPT) are also autoregressive models. They generate text one token at a time, using previous tokens as context.

Comparison of AR, AE, and RNN

Model	AR (Autoregressive)	AE (Autoencoder)	RNN (Recurrent Neural Network)
Goal	Model the joint distribution by predicting each element from past elements	Learn a compressed latent representation and reconstruct the input	Model sequential dependencies using a recurrent hidden state
Input	Past values or past tokens	Entire input vector at once	Input fed one step at a time (sequence)
Memory	Yes (explicit dependence on past tokens)	No explicit memory mechanism	Yes (hidden state stores past information)
Output	Next value/token in the sequence	Reconstructed input	Sequence outputs (one per time step or final state)
Example Use	Language modeling, time series forecasting	Denoising, dimensionality reduction, anomaly detection	Machine translation, speech recognition, sequence classification
Training	Self-supervised (predict next token from context)	Self-supervised / unsupervised (reconstruction loss)	Typically supervised, but can be self-supervised (e.g., language modeling)

Summary of Autoregressive Models

- Sampling: slow (sequential, cannot be parallelized)
 - Sample $\bar{x}_0 \sim p(x_0)$
 - Sample $\bar{x}_1 \sim p(x_1 | x_0 = \bar{x}_0)$
 - Sample $\bar{x}_2 \sim p(x_2 | x_0 = \bar{x}_0, x_1 = \bar{x}_1)$
 - ...
- Density estimation (for training): fast (can be parallelized)
 - Compute $p(x_0 = \bar{x}_0)$
 - Compute $p(x_1 = \bar{x}_1 | x_0 = \bar{x}_0)$
 - Compute $p(x_2 = \bar{x}_2 | x_0 = \bar{x}_0, x_1 = \bar{x}_1)$
 - ...
 - Multiply all terms together (or sum their logarithms), i.e.,

$$\begin{aligned}
 p(\mathbf{x} = \bar{\mathbf{x}}) &= p(x_0 = \bar{x}_0)p(x_1 = \bar{x}_1 | x_0 = \bar{x}_0)p(x_2 = \bar{x}_2 | x_0 = \bar{x}_0, x_1 = \bar{x}_1) \cdots \\
 &\Leftrightarrow \\
 \log p(\mathbf{x} = \bar{\mathbf{x}}) &= \log p(x_0 = \bar{x}_0) + \log p(x_1 = \bar{x}_1 | x_0 = \bar{x}_0) + \log p(x_2 = \bar{x}_2 | x_0 = \bar{x}_0, x_1 = \bar{x}_1) + \cdots
 \end{aligned}$$

- Easy to extend to continuous variables. E.g.,

- Gaussian conditional

$$p(x_t | \mathbf{x}_{<t}) = \mathcal{N}(\mu_\omega(\mathbf{x}_{<t}), \sigma_\omega(\mathbf{x}_{<t}))$$

2. Mixture of logistics

3 Maximum Likelihood Learning (MLE)

3.1 Introduction

How to Understand Generative Model

- Learning a generative model = Density estimation
- Once the full probability distribution is learned, the model can answer any probabilistic inference query.

Given training examples $\{\mathbf{x}^{(i)}\}$ drawn from an unknown true data distribution $p_{\text{data}}(\mathbf{x})$, the goal of a generative model is to learn a parameterized distribution $p_{\theta}(\mathbf{x})$ that approximates $p_{\text{data}}(\mathbf{x})$.

To achieve this, a discrepancy measure $d(p_{\text{data}}(\mathbf{x}), p_{\theta}(\mathbf{x}))$ is minimized. This measures how different the model distribution $p_{\theta}(\mathbf{x})$ is from the true distribution $p_{\text{data}}(\mathbf{x})$.

Common choices for the discrepancy $d(\cdot)$ include:

- Kullback-Leibler (KL) divergence: likelihood-based generative models
- JS divergence: GANs
- f -divergence: f -GANs

However, because $p_{\text{data}}(\mathbf{x})$ is unknown, the discrepancy $d(p_{\text{data}}(\mathbf{x}), p_{\theta}(\mathbf{x}))$ cannot be computed exactly and therefore cannot be used directly as a loss function.

Instead, a practical loss function is defined that satisfies two key properties:

1. It can be computed from available data examples $\{\mathbf{x}^{(i)}\}$.
2. Minimizing the loss encourages the model distribution $p_{\theta}(\mathbf{x})$ to become closer to the true data distribution $p_{\text{data}}(\mathbf{x})$, thereby reducing the underlying discrepancy $d(p_{\text{data}}(\mathbf{x}), p_{\theta}(\mathbf{x}))$.

Relationship Between Loss and KL Divergence

This part explains how the loss function in likelihood-based generative models relates to the KL divergence, thereby justifying its form by showing that the loss satisfies the two desired properties.

For likelihood-based generative models (e.g., autoregressive model, VAEs, normalizing flows), the loss function is defined as the negative log-likelihood of the training data,

$$\mathcal{L}(\theta) = - \sum_{i=1}^n \log p_{\theta}(\mathbf{x}^{(i)})$$

- Minimizing this loss encourages the model distribution p_{θ} to assign high probability to samples drawn from the true data distribution p_{data} .
- Samples \mathbf{x} for which $p_{\theta}(\mathbf{x}) \approx 0$ incur very large penalties because the term $-\log p_{\theta}(\mathbf{x}) \rightarrow \infty$ as $p_{\theta}(\mathbf{x}) \rightarrow 0$.

By calculation,

$$\begin{aligned}
 D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\log \left(\frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right) \right] \\
 &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(p_{\text{data}}(\mathbf{x})) - \log(p_{\theta}(\mathbf{x}))] \\
 &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(p_{\text{data}}(\mathbf{x}))] - \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(p_{\theta}(\mathbf{x}))] \\
 &\approx \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(p_{\text{data}}(\mathbf{x}))] + \mathcal{L}(\theta) \\
 &= -\mathbb{H}(P_{\text{data}}) + \mathcal{L}(\theta)
 \end{aligned}$$

Here,

$$\mathbb{H}(P_{\text{data}}) = -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(p_{\text{data}}(\mathbf{x}))]$$

is the entropy term of the true data distribution p_{data} .

Since the entropy term does not depend on θ ,

$$\begin{aligned}
 \arg \min_{\theta} \mathcal{L}(\theta) &\approx \arg \min_{\theta} (D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) - \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(p_{\text{data}}(\mathbf{x}))]) \\
 &= \arg \min_{\theta} D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x}))
 \end{aligned}$$

However, we cannot directly measure how close we are to the true optimum, since the optimization ignores the entropy term $\mathbb{H}(P_{\text{data}})$. Because this term is unknown and cannot be computed, the absolute value of the KL divergence and therefore the actual distance to the optimum remains unobservable.

Minimizing loss function (i.e., Maximizing Log-likelihood) = Minimizing KL divergence

$$\arg \min_{\theta} \mathcal{L}(\theta) = \arg \max_{\theta} \sum_{i=1}^n \log p_{\theta}(\mathbf{x}^{(i)}) = \arg \min_{\theta} D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x}))$$

KL Divergence

The KL divergence is defined as,

$$D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\log \left(\frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right) \right]$$

The properties of KL divergence:

- Non-negative:

$$\begin{aligned}
 D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\log \left(\frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right) \right] \\
 &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[-\log \left(\frac{p_{\theta}(\mathbf{x})}{p_{\text{data}}(\mathbf{x})} \right) \right] \\
 &\geq \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[1 - \frac{p_{\theta}(\mathbf{x})}{p_{\text{data}}(\mathbf{x})} \right] \\
 &= 1 - \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\frac{p_{\theta}(\mathbf{x})}{p_{\text{data}}(\mathbf{x})} \right] \\
 &= 1 - \int \frac{p_{\theta}(\mathbf{x})}{p_{\text{data}}(\mathbf{x})} \cdot p_{\text{data}}(\mathbf{x}) d\mathbf{x} \\
 &= 1 - \int p_{\theta}(\mathbf{x}) d\mathbf{x} \\
 &= 1 - 1 \\
 &= 0
 \end{aligned}$$

Note: $\ln x \leq x - 1$.

- Non-symmetric: $D_{\text{KL}}(P_{\text{data}} \parallel P_{\theta}) \neq D_{\text{KL}}(P_{\theta} \parallel P_{\text{data}})$
- $D_{\text{KL}}(P_{\text{data}} \parallel P_{\theta}) = 0$ iff the two distributions are the same.
- It measures the compression loss (in bits) of using P_{θ} to encode data that actually comes from P_{data} .

Goal of Learning for Generative Model

Goal of learning

- The goal of learning is to return a model P_{θ} that precisely captures the distribution P_{data} from which our data was sampled
- This is in general not achievable because of
 - limited data only provides a rough approximation of the true underlying distribution
 - computational reasons
- Example. Suppose we represent each image with a vector X of 784 binary variables (black vs. white pixel). How many possible states (= possible images) in the model? $2^{784} \approx 10^{236}$. Even 10^7 training examples provide *extremely sparse coverage*!
- We want to select P_{θ} to construct the "best" approximation to the underlying distribution P_{data}
- What is "best"?

What is “best”?

This depends on what we want to do

- ➊ Density estimation: we are interested in the full distribution (so later we can compute whatever conditional probabilities we want)
- ➋ Specific prediction tasks: we are using the distribution to make a prediction
 - Is this email spam or not?
 - **Structured prediction:** Predict next frame in a video, or caption given an image
- ➌ Structure or knowledge discovery: we are interested in the model itself
 - How do some genes interact with each other?
 - What causes cancer?
 - Take CS 228

Properties of the Monte Carlo Estimate

Properties of the Monte Carlo Estimate

- **Unbiased:**

$$E_P[\hat{g}] = E_P[g(x)]$$

- **Convergence:** By law of large numbers

$$\hat{g} = \frac{1}{T} \sum_{t=1}^T g(x^t) \rightarrow E_P[g(x)] \text{ for } T \rightarrow \infty$$

- **Variance:**

$$V_P[\hat{g}] = V_P \left[\frac{1}{T} \sum_{t=1}^T g(x^t) \right] = \frac{V_P[g(x)]}{T}$$

Thus, variance of the estimator can be reduced by increasing the number of samples.

3.2 MLE for Autoregressive Models

Define the probability function for an autoregressive model with n variables,

$$p_{\theta}(\mathbf{x}) = \prod_{i=1}^n p_{\text{neural}}(x_i | \mathbf{x}_{<i}; \theta_i),$$

with $\theta = (\theta_1, \dots, \theta_n)$ are parameters of all the conditionals.

Let $\mathcal{D} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$. Define the likelihood function,

$$L(\theta, \mathcal{D}) = \prod_{j=1}^m P_{\theta}(\mathbf{x}^{(j)}) = \prod_{j=1}^m \prod_{i=1}^n P_{\text{neutral}}(x_i^{(j)} | \mathbf{x}_{<i}^{(j)}; \theta_i)$$

Define

$$\ell(\theta) = \log L(\theta, \mathcal{D})$$

The optimization turns from $\max_{\theta} L(\theta, \mathcal{D})$ to $\max_{\theta} \ell(\theta)$, that is,

$$\begin{aligned} \arg \max_{\theta} \ell(\theta) &= \arg \max_{\theta} \log L(\theta, \mathcal{D}) \\ &= \arg \max_{\theta} \sum_{j=1}^m \sum_{i=1}^n \log P_{\text{neutral}}(x_i^{(j)} | \mathbf{x}_{<i}^{(j)}; \theta_i) \end{aligned}$$

1. Gradient descent

- (a) Initialize $\theta^0 = (\theta_1, \dots, \theta_n)$
- (b) Compute $\nabla_{\theta} \ell(\theta)$ (by back propagation)
- (c) $\theta^{t+1} = \theta^t + \alpha_t \nabla_{\theta} \ell(\theta)$

2. Stochastic gradient descent

MLE Learning: Stochastic Gradient Descent

$$\ell(\theta) = \log L(\theta, \mathcal{D}) = \sum_{j=1}^m \sum_{i=1}^n \log p_{\text{neutral}}(x_i^{(j)} | \mathbf{x}_{<i}^{(j)}; \theta_i)$$

- ➊ Initialize θ^0 at random
- ➋ Compute $\nabla_{\theta} \ell(\theta)$ (by back propagation)
- ➌ $\theta^{t+1} = \theta^t + \alpha_t \nabla_{\theta} \ell(\theta)$

What is the gradient with respect to θ_i ?

$$\nabla_{\theta_i} \ell(\theta) = \sum_{j=1}^m \nabla_{\theta_i} \sum_{i=1}^n \log p_{\text{neutral}}(x_i^{(j)} | \mathbf{x}_{<i}^{(j)}; \theta_i) = \sum_{j=1}^m \nabla_{\theta_i} \log p_{\text{neutral}}(x_i^{(j)} | \mathbf{x}_{<i}^{(j)}; \theta_i)$$

Each conditional $p_{\text{neutral}}(x_i | \mathbf{x}_{<i}; \theta_i)$ can be optimized separately if there is no parameter sharing. In practice, parameters θ_i are shared (e.g., NADE, PixelRNN, PixelCNN, etc.)

MLE Learning: Stochastic Gradient Descent

$$\ell(\theta) = \log L(\theta, \mathcal{D}) = \sum_{j=1}^m \sum_{i=1}^n \log p_{\text{neural}}(x_i^{(j)} | \mathbf{x}_{<i}^{(j)}; \theta_i)$$

- ➊ Initialize θ^0 at random
- ➋ Compute $\nabla_{\theta} \ell(\theta)$ (by back propagation)
- ➌ $\theta^{t+1} = \theta^t + \alpha_t \nabla_{\theta} \ell(\theta)$

$$\nabla_{\theta} \ell(\theta) = \sum_{j=1}^m \sum_{i=1}^n \nabla_{\theta} \log p_{\text{neural}}(x_i^{(j)} | \mathbf{x}_{<i}^{(j)}; \theta_i)$$

What if $m = |\mathcal{D}|$ is huge?

$$\begin{aligned} \nabla_{\theta} \ell(\theta) &= m \sum_{j=1}^m \frac{1}{m} \sum_{i=1}^n \nabla_{\theta} \log p_{\text{neural}}(x_i^{(j)} | \mathbf{x}_{<i}^{(j)}; \theta_i) \\ &= m E_{x^{(j)} \sim \mathcal{D}} \left[\sum_{i=1}^n \nabla_{\theta} \log p_{\text{neural}}(x_i^{(j)} | \mathbf{x}_{<i}^{(j)}; \theta_i) \right] \end{aligned}$$

Monte Carlo: Sample $x^{(j)} \sim \mathcal{D}; \nabla_{\theta} \ell(\theta) \approx m \sum_{i=1}^n \nabla_{\theta} \log p_{\text{neural}}(x_i^{(j)} | \mathbf{x}_{<i}^{(j)}; \theta_i)$

3.3 Bias and Variance

Bias-Variance trade off

- If the hypothesis space is very limited, it might not be able to represent P_{data} , even with unlimited data
 - This type of limitation is called **bias**, as the learning is limited on how close it can approximate the target distribution
- If we select a highly expressive hypothesis class, we might represent better the data
 - When we have small amount of data, multiple models can fit well, or even better than the true model. Moreover, small perturbations on \mathcal{D} will result in very different estimates
 - This limitation is call the **variance**.

• Stefano Ermon (AI Lab) Deep Generative Models Lecture 4 21 / 25

How to avoid overfitting?

- Hard constraints, e.g. by selecting a less expressive model family:
 - Smaller neural networks with less parameters
 - Weight sharing



- Soft preference for “simpler” models: **Occam Razor**.

- Augment the objective function with **regularization**:

$$\text{objective}(\mathbf{x}, \mathcal{M}) = \text{loss}(\mathbf{x}, \mathcal{M}) + R(\mathcal{M})$$

- Evaluate generalization performance on a held-out validation set

• Stefano Ermon (AI Lab) Deep Generative Models Lecture 4 23 / 25

4 Latent Variable Models (VAEs)

4.1 Latent Variable Models

Motivation

Latent Variable Models: Motivation

- ➊ Only shaded variables x are observed in the data (pixel values)
- ➋ Latent variables z correspond to high level features
 - If z chosen properly, $p(x|z)$ could be much simpler than $p(x)$
 - If we had trained this model, then we could identify features via $p(z|x)$, e.g., $p(EyeColor = Blue|x)$
- ➌ Challenge: Very difficult to specify these conditionals by hand

● Stefano Ermon (AI Lab) Deep Generative Models Lecture 5 5 / 29

Definition: A latent variable model is a probabilistic model, where

1. The data observed x is assumed to be generated from hidden (latent) variables z .
2. These latent variables capture underlying structure or factors in the data.

Basic structure:

$$p(x, z) = p(z)p(x | z)$$

Here, $p(z)$ is prior over latent variables, and $p(x | z)$ is likelihood (how data x is generated from z).

1. Allow us to define complex models $p(x)$ in terms of simpler building blocks $p(x | z)$.
2. Natural for unsupervised learning tasks (clustering, unsupervised representation learning, etc.)
3. No free lunch: much more difficult to learn compared to fully observed, autoregressive models

Example:

1. A Mixture of Gaussians (or Gaussian Mixture Model (GMM)) is a shallow (non-deep, only one level of latent variables) latent-variable model in which a discrete latent variable selects one of K Gaussian components:

$$z \sim \text{Categorical}(\pi_1, \dots, \pi_K)$$

$$x | z = k \sim \mathcal{N}(\mu_k, \Sigma_k)$$

Here, Σ_k is covariance matrix.

The marginal distribution is obtained by summing over the latent variable:

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) = \sum_{k=1}^K p(\mathbf{z} = k)p(\mathbf{x} | \mathbf{z} = k) = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}; \mu_k, \Sigma_k)$$

2. The Variational Autoencoder (VAE) is a type of latent variable model.

4.2 Notations of VAEs

Notation	Explanation
$\mathbf{x} = (x_1, \dots, x_D) \in \mathbb{R}^D$	Observed data (e.g., images, texts, etc.)
$\mathbf{z} = (z_1, \dots, z_L) \in \mathbb{R}^L$	Latent (unobserved) variable representing hidden structure or features
θ	Parameters of the decoder (generative model) <ul style="list-style-type: none"> It is used in the likelihood $p_\theta(\mathbf{x} \mathbf{z})$ and sometimes the prior $p_\theta(\mathbf{z})$ (e.g., when $p_\theta(\mathbf{z}) \sim \mathcal{N}(\mathbf{z} \mathbf{0}, \mathbf{I})$).
ϕ	Parameters of the encoder (inference model) <ul style="list-style-type: none"> It is used in the approximate posterior $q_\phi(\mathbf{z} \mathbf{x})$.
$p_\theta(\mathbf{x})$ or $p(\mathbf{x})$	Marginal likelihood (or true likelihood) of the observed data <ul style="list-style-type: none"> It's called marginal because it is obtained by integrating out (marginalizing over) the latent variable \mathbf{z}: $p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p_\theta(\mathbf{x} \mathbf{z}) p_\theta(\mathbf{z}) d\mathbf{z} = \mathbb{E}_{p_\theta(\mathbf{z})}[p_\theta(\mathbf{x} \mathbf{z})]$ <p>Here, $p_\theta(\mathbf{x}, \mathbf{z})$ is the joint distribution of observed and latent variables, which is also parameterized by θ since</p> $p_\theta(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{x} \mathbf{z}) p_\theta(\mathbf{z})$ <ul style="list-style-type: none"> It is intractable to compute exactly, because \mathbf{z} is high-dimensional.
$\log p_\theta(\mathbf{x})$	Marginal log-likelihood (or true log-likelihood) of the observed data <ul style="list-style-type: none"> It is approximated using variational lower bounds (e.g., ELBO, IWAE bound, etc), since it is intractable to compute exactly.
$p_\theta(\mathbf{x} \mathbf{z})$	Likelihood (or decoder), parameterized by θ <ul style="list-style-type: none"> It models how the observed data \mathbf{x} is generated given the latent variable \mathbf{z}. It typically has a simpler form than the marginal likelihood $p_\theta(\mathbf{x})$. Parameters are obtained from the decoder neural network with parameters θ.
$p_\theta(\mathbf{z})$	Prior distribution over the latent variable, sometimes parameterized by θ
$p_\theta(\mathbf{z} \mathbf{x})$	True posterior distribution of the latent variable \mathbf{z} given the observation \mathbf{x} <ul style="list-style-type: none"> It is defined as $p_\theta(\mathbf{z} \mathbf{x}) = \frac{p_\theta(\mathbf{x}, \mathbf{z})}{p_\theta(\mathbf{x})} = \frac{p_\theta(\mathbf{x} \mathbf{z}) p_\theta(\mathbf{z})}{\int p_\theta(\mathbf{x} \mathbf{z}) p_\theta(\mathbf{z}) d\mathbf{z}}$ <ul style="list-style-type: none"> This posterior characterizes how likely each latent variable \mathbf{z} is, given the observation \mathbf{x}. Although it is the quantity we aim to infer, it is generally intractable and thus approximated by a variational distribution $q_\phi(\mathbf{z} \mathbf{x})$ (as in VAE).
$q_\phi(\mathbf{z} \mathbf{x})$	Approximate posterior (or encoder), parameterized by ϕ <ul style="list-style-type: none"> It approximates the true posterior $p_\theta(\mathbf{z} \mathbf{x})$. Parameters (e.g., mean and log-variance) are obtained from the encoder neural network with parameters ϕ.
$\mathcal{L}_m(\mathbf{x}; \theta, \phi)$	IWAE bound, a tighter lower bound on $\log p_\theta(\mathbf{x})$ that uses m importance-weighted samples
m (sometimes iw)	Number of Monte Carlo samples $\mathbf{z}^{(i)} \sim q_\phi(\mathbf{z} \mathbf{x})$ used per data point to estimate the IWAE bound
$w_i = \frac{p_\theta(\mathbf{x}, \mathbf{z}^{(i)})}{q_\phi(\mathbf{z}^{(i)} \mathbf{x})}$, $i = 1, \dots, m$	Importance weights used to reweight samples in the IWAE objective

4.3 Evidence Lower Bound (ELBO)

Likelihood Intractability

The marginal likelihood is

$$p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x} | \mathbf{z}) p_\theta(\mathbf{z}) d\mathbf{z}$$

However, this integral is generally intractable, because the latent \mathbf{z} is high-dimensional and the decoder $p_\theta(\mathbf{x} | \mathbf{z})$ is nonlinear (e.g., a neural network).

If the true posterior $p_\theta(\mathbf{z} | \mathbf{x})$ were known, the marginal likelihood can be rewritten as an expectation over the posterior,

$$\begin{aligned} p_\theta(\mathbf{x}) &= \int \frac{p_\theta(\mathbf{x} | \mathbf{z}) p_\theta(\mathbf{z})}{p_\theta(\mathbf{z} | \mathbf{x})} \cdot p_\theta(\mathbf{z} | \mathbf{x}) d\mathbf{z} \\ &= \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \mathbf{x})} \left[\frac{p_\theta(\mathbf{x} | \mathbf{z}) p_\theta(\mathbf{z})}{p_\theta(\mathbf{z} | \mathbf{x})} \right] \end{aligned}$$

In principle, this expectation could be approximated via Monte Carlo sampling if able to sample from the true posterior $p_\theta(\mathbf{z} | \mathbf{x})$. However, sampling from the true posterior is itself intractable, because depends on unknown $p_\theta(\mathbf{x})$.

Variational Inference

To address this, a tractable approximate posterior (encoder) $q_\phi(\mathbf{z} | \mathbf{x})$ is introduced. Introducing a tractable approximate posterior is known as variational inference (VI).

By calculation,

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

The second term is a KL divergence between the approximate posterior and the true posterior, which measures how well $q_\phi(\mathbf{z} | \mathbf{x})$ approximates $p_\theta(\mathbf{z} | \mathbf{x})$.

Definition & Property

The Evidence Lower Bound (ELBO) is defined as,

$$\mathcal{L}_{\text{ELBO}}(\mathbf{x}; \theta, \phi) = \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} \left[\log \left(\frac{p_\theta(\mathbf{x} | \mathbf{z}) \cdot p_\theta(\mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x})} \right) \right]$$

ELBO forms a lower bound on the true log-likelihood,

$$\log p_\theta(\mathbf{x}) \geq \mathcal{L}_{\text{ELBO}}(\mathbf{x}; \theta, \phi)$$

In practice, the ELBO is maximized rather than the true log-likelihood, because the true posterior is intractable. The ELBO serves as a tractable, differentiable surrogate objective for learning parameters θ and ϕ . The closer the approximate posterior $q_\phi(\mathbf{z} | \mathbf{x})$ is to the true posterior $p_\theta(\mathbf{z} | \mathbf{x})$, the smaller the KL divergence $D_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}) \| p_\theta(\mathbf{z}))$ becomes, and the tighter the ELBO is as a lower bound to the true log-likelihood.

Decomposition & Calculation

ELBO can be rewritten as

$$\begin{aligned} \mathcal{L}_{\text{ELBO}}(\mathbf{x}; \theta, \phi) &= \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} \left[\log \left(\frac{p_\theta(\mathbf{x} | \mathbf{z}) p_\theta(\mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x})} \right) \right] \\ &= \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} \left[\log p_\theta(\mathbf{x} | \mathbf{z}) - \log \left(\frac{q_\phi(\mathbf{z} | \mathbf{x})}{p_\theta(\mathbf{z})} \right) \right] \\ &= \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} [\log p_\theta(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} \left[\log \left(\frac{q_\phi(\mathbf{z} | \mathbf{x})}{p_\theta(\mathbf{z})} \right) \right] \\ &= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} [\log p_\theta(\mathbf{x} | \mathbf{z})]}_{\text{Reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}) \| p_\theta(\mathbf{z}))}_{\text{KL divergence (regularization) term}} \\ &\approx \underbrace{\frac{1}{S} \sum_{s=1}^S \log p_\theta(\mathbf{x} | \mathbf{z}^{(s)})}_{\text{Monte Carlo approximation}} - D_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}) \| p_\theta(\mathbf{z})) \\ &\approx \log p_\theta(\mathbf{x} | \mathbf{z}^{(1)}) - D_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}) \| p_\theta(\mathbf{z})) \end{aligned}$$

Here, $\mathbf{z}^{(s)} \sim q_\phi(\mathbf{z} | \mathbf{x})$ for $s = 1, \dots, S$.

When $q_\phi(\mathbf{z} | \mathbf{x})$ and $p_\theta(\mathbf{z})$ are both Gaussian distributions,

$$q_\phi(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mu_\phi(\mathbf{x}), \text{diag}(\sigma_\phi^2(\mathbf{x})))$$

$$p_\theta(\mathbf{z}) = \mathcal{N}(0, I)$$

the KL divergence term becomes

$$D_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}) \| p_\theta(\mathbf{z})) = \frac{1}{2} \sum_{i=1}^L (\mu_i^2 + \sigma_i^2 - 1 - \log \sigma_i^2)$$

which is analytically computable (i.e., with a closed-form expression).

Training Objective

By calculation,

$$-\mathcal{L}_{\text{ELBO}}(\mathbf{x}; \theta, \phi) = \underbrace{-\mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} [\log p_\theta(\mathbf{x} | \mathbf{z})]}_{\text{Reconstruction loss}} + \underbrace{D_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}) \| p_\theta(\mathbf{z}))}_{\text{KL divergence term}}$$

For a dataset \mathcal{D} , the learning objective becomes

$$\begin{aligned} \arg \max_{\theta} \sum_{\mathbf{x}^{(i)} \in \mathcal{D}} \log p_{\theta}(\mathbf{x}^{(i)}) &\geq \arg \max_{\theta, \phi^{(1)}, \dots, \phi^{(n)}} \sum_{i=1}^n \mathcal{L}_{\text{ELBO}}(\mathbf{x}^{(i)}; \theta, \phi^{(i)}) \\ &= \arg \min_{\theta, \phi^{(1)}, \dots, \phi^{(n)}} \sum_{i=1}^n -\mathcal{L}_{\text{ELBO}}(\mathbf{x}^{(i)}; \theta, \phi^{(i)}) \end{aligned}$$

Note that use different set of variational parameters $\phi^{(i)}$ for every data point $\mathbf{x}^{(i)}$.

4.4 β -VAE

Definition: β -VAE is a variation of the normal VAE, which is defined as

$$\mathcal{L}_{\beta\text{-ELBO}}(\mathbf{x}; \theta, \phi) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x} | \mathbf{z})] - \beta D_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}) \| p_\theta(\mathbf{z}))$$

Property

In β -VAE, the parameter β controls how much the model focuses on making good reconstructions versus keeping the latent space simple and organized.

1. When $\beta = 1$, β -VAE becomes the normal VAE. The model tries to do two things at once:
 - (a) Reconstruct the input data \mathbf{z} as accurately as possible, i.e., make the reconstruction term as large as possible,
$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x} | \mathbf{z})]$$
 - (b) Make sure the latent variables \mathbf{z} follow the prior distribution $p_\theta(\mathbf{z})$, i.e., make the KL term as small as possible,
$$D_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}) \| p_\theta(\mathbf{z}))$$
2. When $\beta > 1$, the term KL term is multiplied by a larger number. The stronger KL penalty forces the encoder distribution $q_\phi(\mathbf{z} | \mathbf{x})$ to stay close to the prior $p_\theta(\mathbf{z})$, which assumes that all latent dimensions are independent and have simple Gaussian structure. As a result, the latent variables become simpler and more independent, meaning that different dimensions of \mathbf{z} capture different aspects of the data with less overlap between them. This often leads the model to learn more interpretable and disentangled latent features, i.e., each latent variable tends to represent one clear factor of variation in the data (e.g., color, shape, or position). Meanwhile, because the model focuses more on keeping \mathbf{z} better structured (simpler, smoother, more disentangled), it cannot devote as much capacity to maximizing the reconstruction term. As a result, the reconstructed data may become less accurate or blurrier.

To summarize, increasing β produces a latent space that is simpler and more disentangled, but this comes at the cost of reduced reconstruction quality.

4.5 Importance Weighted Autoencoder (IWAE) Bound

Definition: Let $\mathcal{L}_m(\mathbf{x}; \theta, \phi)$ denote the Importance Weighted Autoencoder (IWAE) bound, which is defined as

$$\begin{aligned}\mathcal{L}_m(\mathbf{x}; \theta, \phi) &= \mathbb{E}_{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(\mathbf{z} | \mathbf{x})} \left[\log \left(\frac{1}{m} \sum_{i=1}^m w_i \right) \right] \\ &= \mathbb{E}_{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(\mathbf{z} | \mathbf{x})} \left[\log \left(\frac{1}{m} \sum_{i=1}^m \frac{p_\theta(\mathbf{x}, \mathbf{z}^{(i)})}{q_\phi(\mathbf{z}^{(i)} | \mathbf{x})} \right) \right] \\ &\approx \underbrace{\frac{1}{S} \sum_{s=1}^S \log \left(\frac{1}{m} \sum_{i=1}^m \frac{p_\theta(\mathbf{x}, \mathbf{z}^{(i,s)})}{q_\phi(\mathbf{z}^{(i,s)} | \mathbf{x})} \right)}_{\text{Monte Carlo approximation}} \\ &\approx \log \left(\frac{1}{m} \sum_{i=1}^m \frac{p_\theta(\mathbf{x}, \mathbf{z}^{(i)})}{q_\phi(\mathbf{z}^{(i)} | \mathbf{x})} \right) \\ &= \log \left(\frac{1}{m} \sum_{i=1}^m \frac{p_\theta(\mathbf{x} | \mathbf{z}^{(i)}) p_\theta(\mathbf{z}^{(i)})}{q_\phi(\mathbf{z}^{(i)} | \mathbf{x})} \right)\end{aligned}$$

Here, $\mathbf{z}^{(i,s)} \sim q_\phi(\mathbf{z} | \mathbf{x})$ for $i = 1, \dots, m$ and $s = 1, \dots, S$. m is called the number of importance-weighted samples. w_i (for $i = 1, \dots, m$) denotes the importance weights

$$w_i = \frac{p_\theta(\mathbf{x}, \mathbf{z}^{(i)})}{q_\phi(\mathbf{z}^{(i)} | \mathbf{x})}$$

Property

1. For any $m \geq 1$,

$$\log p_\theta(\mathbf{x}) \geq \mathcal{L}_m(\mathbf{x}; \theta, \phi) \geq \mathcal{L}_1(\mathbf{x}; \theta, \phi) = \mathcal{L}_{\text{ELBO}}(\mathbf{x}; \theta, \phi)$$

Thus, the IWAE bound is always a valid lower bound on the true log-likelihood, and using multiple importance samples ($m > 1$) yields a tighter bound than the ELBO.

- IWAE bound: tighter lower bound obtained with multiple samples.
- ELBO: lower bound corresponding to a single sample.

2. As the number of importance samples grows,

$$\lim_{m \rightarrow \infty} \mathcal{L}_m(\mathbf{x}; \theta, \phi) = \log p_\theta(\mathbf{x})$$

so the IWAE bound converges to the true marginal log-likelihood.

Proof for Property 1

1. Prove that IWAE bound is a lower bound of the true log-likelihood $\log p_\theta(\mathbf{x})$, i.e.,

$$\log p_\theta(\mathbf{x}) \geq \mathcal{L}_m(\mathbf{x}; \theta, \phi)$$

By calculation,

$$\begin{aligned}
\mathbb{E}_{z^{(1)}, \dots, z^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(z|x)} \left[\frac{1}{m} \sum_{i=1}^m w_i \right] &= \mathbb{E}_{z^{(1)}, \dots, z^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(z|x)} \left[\frac{1}{m} \sum_{i=1}^m \frac{p_\theta(x, z^{(i)})}{q_\phi(z^{(i)} | x)} \right] \\
&= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{z^{(i)} \sim q_\phi(z|x)} \left[\frac{p_\theta(x, z^{(i)})}{q_\phi(z^{(i)} | x)} \right] \\
&= \frac{1}{m} \cdot m \cdot \mathbb{E}_{q_\phi(z|x)} \left[\frac{p_\theta(x, z)}{q_\phi(z | x)} \right] \\
&= \mathbb{E}_{q_\phi(z|x)} \left[\frac{p_\theta(x, z)}{q_\phi(z | x)} \right] \\
&= \int \frac{p_\theta(x, z)}{q_\phi(z | x)} q_\phi(z | x) dz \\
&= \int p_\theta(x, z) dz \\
&= p_\theta(x)
\end{aligned}$$

That is, $\frac{1}{m} \sum_{i=1}^m w_i$ is an unbiased estimator of the marginal likelihood $p_\theta(x)$.

Furthermore, according to Jensen's inequality,

$$\begin{aligned}
\log p_\theta(x) &= \log \left(\mathbb{E}_{z^{(1)}, \dots, z^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(z|x)} \left[\frac{1}{m} \sum_{i=1}^m w_i \right] \right) \\
&\geq \mathbb{E}_{z^{(1)}, \dots, z^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(z|x)} \left[\log \left(\frac{1}{m} \sum_{i=1}^m w_i \right) \right] \\
&= \mathcal{L}_m(x; \theta, \phi)
\end{aligned}$$

2. Prove that

$$\mathcal{L}_m(x; \theta, \phi) \geq \mathcal{L}_1(x; \theta, \phi)$$

Let W be a discrete random variable with values in the set $\{w_1, \dots, w_m\}$, each with equal probability $\frac{1}{m}$, i.e., $P(W = w_i) = \frac{1}{m}$, for $i = 1, \dots, m$. By Jensen's inequality,

$$\log \left(\frac{1}{m} \sum_{i=1}^m w_i \right) = \log \mathbb{E}[W] \geq \mathbb{E}[\log W] = \frac{1}{m} \sum_{i=1}^m \log w_i$$

Therefore, by calculation,

$$\begin{aligned}
\mathcal{L}_m(x; \theta, \phi) &= \mathbb{E}_{z^{(1)}, \dots, z^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(z|x)} \left[\log \left(\frac{1}{m} \sum_{i=1}^m w_i \right) \right] \\
&\geq \mathbb{E}_{z^{(1)}, \dots, z^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(z|x)} \left[\frac{1}{m} \sum_{i=1}^m \log w_i \right] \\
&= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{z^{(i)} \sim q_\phi(z|x)} [\log w_i] \\
&= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{z^{(1)} \sim q_\phi(z|x)} [\log w_1] \\
&= \frac{1}{m} \cdot m \cdot \mathbb{E}_{z^{(1)} \sim q_\phi(z|x)} [\log w_1] \\
&= \mathbb{E}_{z^{(1)} \sim q_\phi(z|x)} [\log w_1] \\
&= \mathcal{L}_1(x; \theta, \phi)
\end{aligned}$$

3. Prove that

$$\mathcal{L}_1(\mathbf{x}; \theta, \phi) = \mathcal{L}_{\text{ELBO}}(\mathbf{x}; \theta, \phi)$$

Let $\mathcal{L}_{\text{ELBO}}$ denote the ELBO, by the definition,

$$\begin{aligned} \mathcal{L}_{\text{ELBO}}(\mathbf{x}; \theta, \phi) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \left(\frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x})} \right) \right] \\ &= \mathbb{E}_{\mathbf{z}^{(1)} \sim q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \left(\frac{p_\theta(\mathbf{x}, \mathbf{z}^{(1)})}{q_\phi(\mathbf{z}^{(1)} | \mathbf{x})} \right) \right] \\ &= \mathbb{E}_{\mathbf{z}^{(1)} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log w_1] \\ &= \mathcal{L}_1(\mathbf{x}; \theta, \phi) \end{aligned}$$

Proof for Property 2

$$\begin{aligned} \lim_{m \rightarrow \infty} \mathcal{L}_m(\mathbf{x}; \theta, \phi) &= \mathbb{E}_{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \left(\frac{1}{m} \sum_{i=1}^m \frac{p_\theta(\mathbf{x}, \mathbf{z}^{(i)})}{q_\phi(\mathbf{z}^{(i)} | \mathbf{x})} \right) \right] \\ &= \mathbb{E}_{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \left(\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} \left[\frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x})} \right] \right) \right] \\ &= \mathbb{E}_{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \left(\int \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x})} q_\phi(\mathbf{z} | \mathbf{x}) d\mathbf{z} \right) \right] \\ &= \mathbb{E}_{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \left(\int p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z} \right) \right] \\ &= \mathbb{E}_{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)} \stackrel{\text{i.i.d.}}{\sim} q_\phi(\mathbf{z}|\mathbf{x})} [\log(p_\theta(\mathbf{x}))] \\ &= \log p_\theta(\mathbf{x}) \end{aligned}$$

4.6 Assumptions of VAEs

	Variational Auto Encoder (VAE)	Mixture of Gaussians VAE (GMVAE)	Importance Weight Auto Encoder (IWAE)
Optimization	ELBO	ELBO	IWAE bound (IW-ELBO)
Prior	$p(z) = \mathcal{N}(z 0, I)$	$p_\theta(z) = \sum_{i=1}^k \frac{1}{k} \mathcal{N}(z \mu_i, \text{diag}(\sigma_i^2))$	$p(z) = \mathcal{N}(z 0, I)$
Likelihood (decoder)		<p>1. Binary data ($x_i \in \{0, 1\}$): Bernoulli distribution $p_\theta(x z) = \text{Bern}(x f_\theta(z))$</p> <p>2. Continuous data ($x_i \in \mathbb{R}$): Gaussian distribution $p_\theta(x z) = \mathcal{N}(x \mu_\theta(z), \Sigma_\theta(z)) = \mathcal{N}(x \mu_\theta(z), I)$</p> <p>Here</p> $\mu_\theta(z) = \sigma(Az + c) \in \mathbb{R}^L$ $\Sigma_\theta(z) = \text{diag}\left(e^{\sigma(Bz+d)}\right) \in \mathbb{R}^L$ $\theta = (A, B, c, d)$ <p>3. Categorical data ($x_i \in \{1, 2, \dots, K\}$): Categorical distribution $p_\theta(x z) = \text{Categorical}(x f_\theta(z))$</p>	
Approximate posterior (encoder)	$q_\phi(z x) = \mathcal{N}(z \mu_\phi(x), \text{diag}(\sigma_\phi^2(x)))$		

4.7 Workflow of VAEs

1. Feed the input, an observation $\mathbf{x} \in \mathbb{R}^D$, into the encoder neural network, which outputs two vectors for each input
 - (a) The mean vector of the encoder $q_\phi(\mathbf{z} | \mathbf{x})$: $\mu_\phi(\mathbf{x}) \in \mathbb{R}^L$
 - (b) The log-variance vector of the encoder $q_\phi(\mathbf{z} | \mathbf{x})$: $\log \sigma_\phi^2(\mathbf{x}) \in \mathbb{R}^L$

2. Compute the encoder $q_\phi(\mathbf{z} | \mathbf{x})$, i.e.,

$$q_\phi(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \mu_\phi(\mathbf{x}), \text{diag}(\sigma_\phi^2(\mathbf{x})))$$

3. Sample $\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x})$, by reparameterization trick.

It is impossible to directly backpropagate through the sampling process

$$\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \mu_\phi(\mathbf{x}), \text{diag}(\sigma_\phi^2(\mathbf{x})))$$

since this stochastic operation is non-differentiable with respect to the parameters ϕ .

To make gradient-based optimization possible, the random variable \mathbf{z} is rewritten as a deterministic function of a parameter-free random variable ϵ and differentiable parameters $\mu_\phi(\mathbf{x})$ and $\sigma_\phi(\mathbf{x})$

$$\mathbf{z} = \mu_\phi(\mathbf{x}) + \sigma_\phi(\mathbf{x}) \odot \epsilon$$

$$\epsilon \sim \mathcal{N}(0, I)$$

Here, \odot represents element-wise product. The randomness is entirely captured by ϵ , which does not depend on ϕ , enabling backpropagation through \mathbf{z} during training.

- (a) VAE: One sample $\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x})$
- (b) IWAE: m samples $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)} \sim q_\phi(\mathbf{z} | \mathbf{x})$

4. Feed $\mathbf{z} \in \mathbb{R}^L$ into the decoder neural network $f_\theta(\cdot)$, which output

- (a) Continuous, $x_i \in \mathbb{R}$: The mean vector of the Gaussian likelihood

$$\mu_\theta(\mathbf{z}) = (\mu_1, \dots, \mu_D) = f_\theta(\mathbf{z}) \in \mathbb{R}^D$$

- (b) Binary (discrete), $x_i \in \{0, 1\}$: The vector of D logits (representing the raw, unnormalized activations), i.e., $(\alpha_1, \dots, \alpha_D)$

$$(\alpha_1, \dots, \alpha_D) = f_\theta(\mathbf{z}) \in \mathbb{R}^D$$

- (c) Categorical (discrete), $x_i \in \{1, 2, \dots, K\}$: The matrix of $D \times K$ logits, i.e.,

$$\ell = f_\theta(\mathbf{z}) \in \mathbb{R}^{D \times K}$$

5. Compute the decoder $p_\theta(\mathbf{x} | \mathbf{z})$. Specifically,

- (a) Continuous, $x_i \in \mathbb{R}$:

$$p_\theta(\mathbf{x} | \mathbf{z}) = \prod_{i=1}^D \mathcal{N}(x_i | \mu_i, 1)$$

- (b) Binary, $x_i \in \{0, 1\}$:

$$p_\theta(\mathbf{x} | \mathbf{z}) = \prod_{i=1}^D \text{Bern}(x_i | p_i)$$

Here, p_i is obtained by applying the sigmoid to logits,

$$p_i = \sigma(\alpha_i) = \frac{1}{1 + e^{-\alpha_i}}$$

(c) Categorical, $x_i \in \{1, 2, \dots, K\}$:

$$p_\theta(\mathbf{x} | \mathbf{z}) = \prod_{i=1}^D \text{Categorical}(x_i | \mathbf{p}_i) = \prod_{i=1}^D \text{Categorical}(x_i | (p_{i,1}, \dots, p_{i,K}))$$

Here, $\mathbf{p}_i = (p_{i,1}, \dots, p_{i,K})$ is obtained by applying the softmax to logits,

$$\mathbf{p}_i = (p_{i,1}, \dots, p_{i,K}) = \left(\frac{e^{\ell_{i,1}}}{\sum_{j=1}^K e^{\ell_{i,j}}}, \dots, \frac{e^{\ell_{i,K}}}{\sum_{j=1}^K e^{\ell_{i,j}}} \right)$$

6. Apply $p_\theta(\mathbf{x} | \mathbf{z})$ to

(a) Approximate (using Mont Carlo) ELBO reconstruction term

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x} | \mathbf{z})] \approx \log p_\theta(\mathbf{x} | \mathbf{z})$$

(b) Approximate (using Mont Carlo) IWAE bound

$$\mathcal{L}_m(\mathbf{x}; \theta, \phi) \approx \log \left(\frac{1}{m} \sum_{i=1}^m \frac{p_\theta(\mathbf{x} | \mathbf{z}^{(i)}) p_\theta(\mathbf{z}^{(i)})}{q_\phi(\mathbf{z}^{(i)} | \mathbf{x})} \right)$$

7. Repeat the above process for n observations \mathbf{x} ,

(a) ELBO

$$(\theta^*, \phi^*) = \arg \max_{\theta, \phi} \sum_{i=1}^n \mathcal{L}_{\text{ELBO}}(\mathbf{x}^{(i)}; \theta, \phi)$$

(b) IWAE bound

$$(\theta^*, \phi^*) = \arg \max_{\theta, \phi} \sum_{i=1}^n \mathcal{L}_m(\mathbf{x}^{(i)}; \theta, \phi)$$

8. Sampling after model training

(a) Sample a latent from the prior, i.e., $\mathbf{z} \sim p_\theta(\mathbf{z})$.

(b) Compute the parameters (mean and variance) of the decoder $p_\theta(\mathbf{x} | \mathbf{z})$ from \mathbf{z} .

(c) Sample an observation from the decoder. i.e., $\mathbf{x} \sim p_\theta(\mathbf{x} | \mathbf{z})$.

5 Normalizing Flows

5.1 Introduction

Change of Variables formula

1. 1-dimensional continuous case

- (a) Let $X, Z \in \mathbb{R}$ be two random variables with probability densities $p_X(x)$ and $p_Z(z)$. Suppose $X = f(Z)$, and $f(\cdot)$ is monotone with inverse $Z = h(X) = f^{-1}(X)$. Then

$$p_X(x) = p_Z(h(x)) |h'(x)| = p_Z(z) \frac{1}{|f'(x)|}$$

2. n -dimensional continuous case, linear transformation

- (a) Let $X, Z \in \mathbb{R}^n$ be two random variables with probability densities $p_X(x)$ and $p_Z(z)$. Suppose $X = AZ$, where A is an invertible square matrix with inverse $W = A^{-1}$. Then

$$p_X(x) = p_Z(Wx) |\det(W)| = p_Z(z) \frac{1}{|\det(A)|}$$

3. n -dimensional continuous case, general transformation

- (a) Let $X, Z \in \mathbb{R}^n$ be two random variables with probability densities $p_X(x)$ and $p_Z(z)$. Suppose $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is invertible such that $X = \mathbf{f}(Z)$ and $Z = \mathbf{f}^{-1}(X)$. Then

$$p_X(x) = p_Z(\mathbf{f}^{-1}(x)) \left| \det \left(\frac{\partial \mathbf{f}^{-1}(x)}{\partial x} \right) \right| = p_Z(z) \left| \det \left(\frac{\partial \mathbf{f}(z)}{\partial z} \right) \right|^{-1}$$

5.2 Comparison of VAE & NF

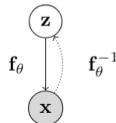
Both Variational Autoencoders (VAE) and Normalizing Flows (NF) consist of an encoder and a decoder. The key distinction is whether the decoder is invertible, which fundamentally impacts the design of the encoder and the computation of the data likelihood.

	VAE	NF
Decoder Definition	The decoder is defined as a conditional probability distribution $p_\theta(x z)$. Typically, a neural network $f_\theta(z)$ outputs the parameters (e.g., mean and variance) of this distribution.	The decoder is defined as a deterministic and invertible transformation, $f_\theta : z \rightarrow x, \quad x = f_\theta(z)$
Decoder Feature	The decoder $p_\theta(x z)$ is stochastic and therefore not invertible.	The decoder f_θ is deterministic and invertible.
Decoder Feature	Given a latent variable z , the output x is stochastic, i.e., multiple possible x samples can be drawn from $p_\theta(x z)$.	Given a latent variable z , the output x is deterministic, i.e., one z corresponds to exactly one x .
Encoder Definition	Given an observation x , there may be multiple possible z values. Therefore, the encoder is defined probabilistically, $z \sim q_\phi(z x)$	Given an observation x , there is exactly one corresponding z . The encoder is defined deterministically, $z = f_\theta^{-1}(x)$
Mapping	There is no one-to-one mapping between x and z .	There is a one-to-one mapping between x and z .
Type of x	Can be continuous or discrete (binary, categorical)	Must be continuous
Type of z	Usually continuous	Must be continuous
Latent Dimension	Since the decoder is not invertible, the latent dimension can differ from the input dimension, allowing a lower-dimensional representation.	Invertibility mathematically requires equal dimensionality. The latent dimension must equal the input dimension, limiting the possibility of a lower-dimensional representation.
Likelihood Computation	The lack of invertibility prevents direct computation of the likelihood $p_\theta(x)$. Instead, a variational lower bound (e.g., ELBO, IWAE bound, etc) is optimized.	Invertibility allows direct computation of the likelihood $p_\theta(x)$, using the change-of-variables formula.
Sampling (Post-training)	<ol style="list-style-type: none"> 1. Sample from the prior $z \sim p_\theta(z)$. 2. Compute the parameters (mean and variance) of the decoder $p_\theta(x z)$, using z. 3. Sample from the decoder. i.e., $x \sim p_\theta(x z)$. 	<ol style="list-style-type: none"> 1. Sample from the prior $z \sim p_\theta(z)$. 2. Compute $x = f_\theta(z)$.
Workflow	$\begin{array}{c} x \xrightarrow[\text{NN: } \mu_\phi(x), \sigma_\phi(x)]{\text{Encoder}} q_\phi(z x) \xrightarrow[z \sim q_\phi(z x)]{\text{Sampling}} z \xrightarrow[\text{NN } f_\theta(z)]{\text{Decoder}} x \\ p_\theta(x z) \xrightarrow{\text{Compute}} \text{Lower bound of } p_\theta(x) \end{array}$	$\begin{array}{c} p_\theta(z) \xrightarrow[z \sim p_\theta(z)]{\text{Sampling}} z \xrightarrow[x=f_\theta(z)]{\text{Decoder}} x \xrightarrow[\text{Change-of-variable}]{\text{Compute}} p_\theta(x) \end{array}$

5.3 Normalizing Flows

Normalizing flow models

- Consider a directed, latent-variable model over observed variables X and latent variables Z
- In a **normalizing flow model**, the mapping between Z and X , given by $\mathbf{f}_\theta : \mathbb{R}^n \mapsto \mathbb{R}^n$, is deterministic and invertible such that $X = \mathbf{f}_\theta(Z)$ and $Z = \mathbf{f}_\theta^{-1}(X)$



- Using change of variables, the marginal likelihood $p(x)$ is given by

$$p_X(x; \theta) = p_Z(\mathbf{f}_\theta^{-1}(x)) \left| \det \left(\frac{\partial \mathbf{f}_\theta^{-1}(x)}{\partial x} \right) \right|$$

- Note: x, z need to be continuous and have the same dimension.

- Stefano Ermon (AI Lab) Deep Generative Models Lecture 7 13 / 19

A Flow of Transformations

Normalizing: Change of variables gives a normalized density after applying an invertible transformation

Flow: Invertible transformations can be composed with each other

$$\mathbf{z}_m = \mathbf{f}_\theta^m \circ \dots \circ \mathbf{f}_\theta^1(\mathbf{z}_0) = \mathbf{f}_\theta^m(\mathbf{f}_\theta^{m-1}(\dots(\mathbf{f}_\theta^1(\mathbf{z}_0)))) \triangleq \mathbf{f}_\theta(\mathbf{z}_0)$$

- Start with a simple distribution for \mathbf{z}_0 (e.g., Gaussian)
- Apply a sequence of M invertible transformations to finally obtain $\mathbf{x} = \mathbf{z}_M$
- By change of variables

$$p_X(x; \theta) = p_Z(\mathbf{f}_\theta^{-1}(x)) \prod_{m=1}^M \left| \det \left(\frac{\partial (\mathbf{f}_\theta^m)^{-1}(\mathbf{z}_m)}{\partial \mathbf{z}_m} \right) \right|$$

(Note: determinant of product equals product of determinants)

- Stefano Ermon (AI Lab) Deep Generative Models Lecture 7 14 / 19

Learning and Inference

- Learning via **maximum likelihood** over the dataset \mathcal{D}

$$\max_{\theta} \log p_X(\mathcal{D}; \theta) = \sum_{x \in \mathcal{D}} \log p_Z(\mathbf{f}_\theta^{-1}(x)) + \log \left| \det \left(\frac{\partial \mathbf{f}_\theta^{-1}(x)}{\partial x} \right) \right|$$

- Exact likelihood evaluation** via inverse transformation $x \mapsto z$ and change of variables formula
- Sampling** via forward transformation $z \mapsto x$

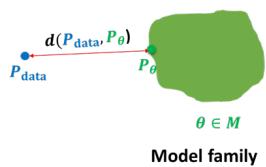
$$z \sim p_Z(z) \quad x = \mathbf{f}_\theta(z)$$

- Latent representations** inferred via inverse transformation (no inference network required!)

$$z = \mathbf{f}_\theta^{-1}(x)$$

- Stefano Ermon (AI Lab) Deep Generative Models Lecture 7 16 / 19

Recap



- Model families
 - Autoregressive Models: $p_\theta(\mathbf{x}) = \prod_{i=1}^n p_\theta(x_i | \mathbf{x}_{<i})$
 - Variational Autoencoders: $p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z}$
 - Normalizing Flow Models: $p_X(\mathbf{x}; \theta) = p_Z(\mathbf{f}_\theta^{-1}(\mathbf{x})) \left| \det \left(\frac{\partial \mathbf{f}_\theta^{-1}(\mathbf{x})}{\partial \mathbf{x}} \right) \right|$
- All the above families are trained by minimizing KL divergence $d_{KL}(p_{\text{data}} \| p_\theta)$, or equivalently maximizing likelihoods (or approximations)
- Today: alternative choices for $d(p_{\text{data}} \| p_\theta)$

6 Generative Adversarial Networks (GANs)

6.1 Motivation of Likelihood-Free Learning

Types of learning

1. Density estimation (maximize likelihood)

- Classical generative models (e.g., AR, VAE, NF) are typically trained by maximizing the log-likelihood $\max_{P_\theta} \sum_{i=1}^n \log P_\theta(\mathbf{x}^{(i)})$, which is equivalent to minimizing the KL divergence between the data and model distributions $\min_{P_\theta} D_{\text{KL}}(P_{\text{data}} \parallel P_\theta)$.

2. Likelihood-free learning (train a binary classifier, which is called discriminator in GAN)

- However, a high log-likelihood (or small KL divergence) does not necessarily imply good sample quality, and vice versa. Therefore, it is important to disentangle likelihood from sample quality, motivating the development of likelihood-free learning approaches.

6.2 Notation for GAN

Notation	Explanation
$D_\phi : \mathbb{R}^D \rightarrow [0, 1]$	<p>Discriminator:</p> <ul style="list-style-type: none"> • It tries to tell apart real samples $x \sim p_{\text{data}}(x)$ from fake samples $x = G_\theta(z)$ with $z \sim p(z)$. • Input: data sample x (real or fake) • Output: probability being real • Typically implemented as a neural network with parameters ϕ, e.g., $D_\phi(x) = \sigma(h_\phi(x))$, <ul style="list-style-type: none"> – Logits $h_\phi(x)$: $x \rightarrow \text{Flatten} \rightarrow \text{Linear} \rightarrow \text{LeakyReLU} \rightarrow \text{Linear}$ – Sigmoid activation $\sigma(\cdot)$: convert logits $h_\phi(x)$ to probability $D_\phi(x)$
$G_\theta(z) : \mathbb{R}^L \rightarrow \mathbb{R}^D$	<p>Generator :</p> <ul style="list-style-type: none"> • It tries to fool the discriminator by generating fake samples that appear real, i.e., $D_\phi(G_\theta(z)) \rightarrow 1$. • Input: latent vector $z \sim p(z)$ • Output: fake sample x • Typically implemented as a neural network with parameters θ.

6.3 GAN

- A Generative Adversarial Network (GAN) consists of two players:
 1. Discriminator $D_\phi(\mathbf{x})$
 2. Generator $G_\theta(\mathbf{z})$
- Each player minimizes its own loss while maximizing against the other.

6.4 Discriminator Loss

- Intuitively, $D_\phi(\mathbf{x})$ should assign
 1. high probability to real samples, $D_\phi(\mathbf{x}) \rightarrow 1$
 2. low probability to fake samples, $D_\phi(G_\theta(\mathbf{z})) \rightarrow 0$
- The loss function is defined as,

$$\begin{aligned}\mathcal{L}_D(\phi; \theta) &= -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D_\phi(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} [\log(1 - D_\phi(\mathbf{x}))] \\ &= -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D_\phi(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D_\phi(G_\theta(\mathbf{z})))]\end{aligned}$$

Here, $p_\theta(\mathbf{x})$ denotes the distribution of fake samples generated by $G_\theta(\mathbf{z})$ with $\mathbf{z} \sim p(\mathbf{z})$ and $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, I)$.

- The goal of the discriminator is to minimize the loss, which is equivalent to maximizing the negative loss,

$$\max_{\phi} -\mathcal{L}_D(\phi; \theta) = \max_{\phi} (\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D_\phi(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D_\phi(G_\theta(\mathbf{z})))])$$

- $\mathcal{L}_D(\phi; \theta)$ is minimized when $D_\phi(\mathbf{x}) = D_\phi^*(\mathbf{x})$, with

$$D_\phi^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}$$

- According to definition of $\mathcal{L}_D(\phi; \theta)$,

$$\begin{aligned}\mathcal{L}_D(\phi; \theta) &= -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D_\phi(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} [\log(1 - D_\phi(\mathbf{x}))] \\ &= - \int p_{\text{data}}(\mathbf{x}) \log D_\phi(\mathbf{x}) d\mathbf{x} - \int p_\theta(\mathbf{x}) \log(1 - D_\phi(\mathbf{x})) d\mathbf{x} \\ &= \int (-p_{\text{data}}(\mathbf{x}) \log D_\phi(\mathbf{x}) - p_\theta(\mathbf{x}) \log(1 - D_\phi(\mathbf{x}))) d\mathbf{x}\end{aligned}$$

Thus, by calculation,

$$\begin{aligned}\mathcal{L}'_D(\phi; \theta) &= \frac{d}{dD_\phi(\mathbf{x})} \int (-p_{\text{data}}(\mathbf{x}) \log D_\phi(\mathbf{x}) - p_\theta(\mathbf{x}) \log(1 - D_\phi(\mathbf{x}))) d\mathbf{x} \\ &= \int \frac{d}{dD_\phi(\mathbf{x})} (-p_{\text{data}}(\mathbf{x}) \log D_\phi(\mathbf{x}) - p_\theta(\mathbf{x}) \log(1 - D_\phi(\mathbf{x}))) d\mathbf{x} \\ &= \int \frac{d}{dt} (-p_{\text{data}}(\mathbf{x}) \log t - p_\theta(\mathbf{x}) \log(1 - t)) \Big|_{t=D_\phi(\mathbf{x})} d\mathbf{x} \\ &= \int \left(\frac{p_\theta(\mathbf{x})}{1-t} - \frac{p_{\text{data}}(\mathbf{x})}{t} \right) \Big|_{t=D_\phi(\mathbf{x})} d\mathbf{x} \\ &= \int \left(\frac{p_\theta(\mathbf{x})}{1-D_\phi(\mathbf{x})} - \frac{p_{\text{data}}(\mathbf{x})}{D_\phi(\mathbf{x})} \right) d\mathbf{x}\end{aligned}$$

and

$$\begin{aligned}
\mathcal{L}_D''(\phi; \theta) &= \frac{d}{dD_\phi(\mathbf{x})} \int \left(\frac{p_\theta(\mathbf{x})}{1 - D_\phi(\mathbf{x})} - \frac{p_{\text{data}}(\mathbf{x})}{D_\phi(\mathbf{x})} \right) d\mathbf{x} \\
&= \int \frac{d}{dD_\phi(\mathbf{x})} \left(\frac{p_\theta(\mathbf{x})}{1 - D_\phi(\mathbf{x})} - \frac{p_{\text{data}}(\mathbf{x})}{D_\phi(\mathbf{x})} \right) d\mathbf{x} \\
&= \int \frac{d}{dt} \left(\frac{p_\theta(\mathbf{x})}{1 - t} - \frac{p_{\text{data}}(\mathbf{x})}{t} \right) \Big|_{t=D_\phi(\mathbf{x})} d\mathbf{x} \\
&= \int \left(\frac{p_\theta(\mathbf{x})}{(1-t)^2} + \frac{p_{\text{data}}(\mathbf{x})}{t^2} \right) \Big|_{t=D_\phi(\mathbf{x})} d\mathbf{x} \\
&= \int \left(\frac{p_\theta(\mathbf{x})}{(1-D_\phi(\mathbf{x}))^2} + \frac{p_{\text{data}}(\mathbf{x})}{D_\phi^2(\mathbf{x})} \right) d\mathbf{x} \\
&\geq 0
\end{aligned}$$

Therefore,

$$\begin{aligned}
\mathcal{L}'_D(\phi; \theta) = 0 &\Rightarrow \int \left(\frac{p_\theta(\mathbf{x})}{1 - D_\phi(\mathbf{x})} - \frac{p_{\text{data}}(\mathbf{x})}{D_\phi(\mathbf{x})} \right) d\mathbf{x} = 0 \\
&\Rightarrow \frac{p_\theta(\mathbf{x})}{1 - D_\phi(\mathbf{x})} - \frac{p_{\text{data}}(\mathbf{x})}{D_\phi(\mathbf{x})} = 0 \\
&\Rightarrow D_\phi(\mathbf{x}) \cdot p_\theta(\mathbf{x}) - (1 - D_\phi(\mathbf{x})) \cdot p_{\text{data}}(\mathbf{x}) = 0 \\
&\Rightarrow D_\phi(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}
\end{aligned}$$

- If $D_\phi(\mathbf{x}) = D_\phi^*(\mathbf{x})$, then

$$h_\phi(\mathbf{x}) = \log \left(\frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x})} \right)$$

– Since

$$\begin{aligned}
D_\phi^*(\mathbf{x}) &= \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})} \\
D_\phi(\mathbf{x}) &= \sigma(h_\phi(\mathbf{x})) = \frac{1}{1 + e^{-h_\phi(\mathbf{x})}}
\end{aligned}$$

Thus, when $D_\phi(\mathbf{x}) = D_\phi^*(\mathbf{x})$,

$$\begin{aligned}
\frac{1}{1 + e^{-h_\phi(\mathbf{x})}} &= \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})} \Rightarrow 1 + e^{-h_\phi(\mathbf{x})} = \frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x})} \\
&\Rightarrow e^{-h_\phi(\mathbf{x})} = \frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x})} - 1 \\
&\Rightarrow e^{-h_\phi(\mathbf{x})} = \frac{p_\theta(\mathbf{x})}{p_{\text{data}}(\mathbf{x})} \\
&\Rightarrow -h_\phi(\mathbf{x}) = \log \left(\frac{p_\theta(\mathbf{x})}{p_{\text{data}}(\mathbf{x})} \right) \\
&\Rightarrow h_\phi(\mathbf{x}) = \log \left(\frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x})} \right)
\end{aligned}$$

6.5 Generator Minimax Loss

- The generator aims to fool the discriminator, i.e., minimizes the same value function that the discriminator tries to maximize,

$$\begin{aligned} & \arg \max_{\theta} (-\mathbb{E}_{x \sim p_{\text{data}}(x)}[\log D_{\phi}(x)] - \mathbb{E}_{z \sim p(z)}[\log(1 - D_{\phi}(G_{\theta}(z)))] \\ &= \arg \min_{\theta} (\mathbb{E}_{x \sim p_{\text{data}}(x)}[\log D_{\phi}(x)] + \mathbb{E}_{z \sim p(z)}[\log(1 - D_{\phi}(G_{\theta}(z)))] \\ &= \arg \min_{\theta} \mathbb{E}_{z \sim p(z)}[\log(1 - D_{\phi}(G_{\theta}(z)))] \end{aligned}$$

- The generator minimax loss is defined as,

$$\begin{aligned} \mathcal{L}_G^{\text{minimax}}(\theta; \phi) &= \mathbb{E}_{x \sim p_{\theta}(x)}[\log(1 - D_{\phi}(x))] \\ &= \mathbb{E}_{z \sim p(z)}[\log(1 - D_{\phi}(G_{\theta}(z)))] \end{aligned}$$

- The goal of the generator is to minimize this loss,

$$\min_{\theta} \mathcal{L}_G^{\text{minimax}}(\theta; \phi) = \min_{\theta} \mathbb{E}_{z \sim p(z)}[\log(1 - D_{\phi}(G_{\theta}(z)))]$$

6.6 Generator Non-Saturating Loss

- The generator non-saturating loss is defined,

$$\begin{aligned}\mathcal{L}_G^{\text{non-sat}}(\theta; \phi) &= -\mathbb{E}_{x \sim p_\theta(x)}[\log D_\phi(x)] \\ &= -\mathbb{E}_{z \sim p(z)}[\log D_\phi(G_\theta(z))]\end{aligned}$$

- The goal of the generator is to minimize this loss, which is equivalent to maximize its negative,

$$\max_{\theta} -\mathcal{L}_G^{\text{non-sat}}(\theta; \phi) = \max_{\theta} \mathbb{E}_{z \sim p(z)}[\log D_\phi(G_\theta(z))]$$

6.7 Generator KL-based Loss

- The generator KL-based loss is defined as the sum of the minimax loss and the non-saturating loss,

$$\begin{aligned}\mathcal{L}_G(\theta; \phi) &= \mathcal{L}_G^{\text{minimax}}(\theta; \phi) + \mathcal{L}_G^{\text{non-sat}}(\theta; \phi) \\ &= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} [\log(1 - D_\phi(\mathbf{x}))] - \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} [\log D_\phi(\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D_\phi(G_\theta(\mathbf{z})))] - \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log D_\phi(G_\theta(\mathbf{z}))]\end{aligned}$$

- If $D_\phi(\mathbf{x}) = D_\phi^*(\mathbf{x})$, then

$$\mathcal{L}_G(\theta; \phi) = D_{\text{KL}}(p_\theta(\mathbf{x}) \parallel p_{\text{data}}(\mathbf{x}))$$

– Since

$$D_\phi^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}$$

Thus, by calculation,

$$\begin{aligned}\mathcal{L}_G(\theta; \phi) &= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} [\log(1 - D_\phi(\mathbf{x}))] - \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} [\log D_\phi(\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \left(\frac{1 - D_\phi(\mathbf{x})}{D_\phi(\mathbf{x})} \right) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \left(\frac{1 - D_\phi^*(\mathbf{x})}{D_\phi^*(\mathbf{x})} \right) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \left(\frac{1 - \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}}{\frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}} \right) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \left(\frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x}) - p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x})} \right) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \left(\frac{p_\theta(\mathbf{x})}{p_{\text{data}}(\mathbf{x})} \right) \right] \\ &= D_{\text{KL}}(p_\theta(\mathbf{x}) \parallel p_{\text{data}}(\mathbf{x}))\end{aligned}$$

6.8 WGAN Critic Loss

- The WGAN critic loss is defined as,

$$\begin{aligned}\mathcal{L}_D(\phi; \theta) &= -\mathbb{E}_{x \sim p_{\text{data}}(x)}[D_\phi(x)] + \mathbb{E}_{x \sim p_\theta(x)}[D_\phi(x)] \\ &= -\mathbb{E}_{x \sim p_{\text{data}}(x)}[D_\phi(x)] + \mathbb{E}_{z \sim p(z)}[D_\phi(G_\theta(z))]\end{aligned}$$

6.9 WGAN-GP Critic Loss

- The WGAN-GP critic loss is defined as,

$$\begin{aligned}\mathcal{L}_D(\phi; \theta) &= -\mathbb{E}_{x \sim p_{\text{data}}(x)}[D_\phi(x)] + \mathbb{E}_{x \sim p_\theta(x)}[D_\phi(x)] + \lambda \mathbb{E}_{x \sim r_\theta(x)}[(\|\nabla D_\phi(x)\|_2 - 1)^2] \\ &= -\mathbb{E}_{x \sim p_{\text{data}}(x)}[D_\phi(x)] + \mathbb{E}_{z \sim p(z)}[D_\phi(G_\theta(z))] + \lambda \mathbb{E}_{x \sim r_\theta(x)}[(\|\nabla D_\phi(x)\|_2 - 1)^2]\end{aligned}$$

6.10 WGAN Generator Loss

- The WGAN generator loss is defined as,

$$\begin{aligned}\mathcal{L}_G(\theta; \phi) &= -\mathbb{E}_{z \sim p_\theta(z)}[D_\phi(x)] \\ &= -\mathbb{E}_{z \sim p(z)}[D_\phi(G_\theta(z))]\end{aligned}$$

6.11 Standard GAN

- Definition:

1. Discriminator loss

- Optimal discriminator,

$$D_\phi^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}$$

2. Generator minimax loss

- Note: Discriminator loss and generator minimax loss can be combined as,

$$\begin{aligned} V(G, D) &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D_\phi(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} [\log(1 - D_\phi(\mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D_\phi(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D_\phi(G_\theta(\mathbf{z})))] \end{aligned}$$

with the training objective

$$\min_G \max_D V(G, D) = \min_G V(G, D^*)$$

By calculation,

$$\begin{aligned} V(G, D^*) &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D_\phi^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} [\log(1 - D_\phi^*(\mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \left[\log \left(\frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})} \right) \right] + \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \left(1 - \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})} \right) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \left[\log \left(\frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})} \right) \right] + \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \left(\frac{p_\theta(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})} \right) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \left[\log \left(\frac{p_{\text{data}}(\mathbf{x})}{\frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2}} \cdot \frac{1}{2} \right) \right] + \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \left(\frac{p_\theta(\mathbf{x})}{\frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2}} \cdot \frac{1}{2} \right) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \left[\log \left(\frac{p_{\text{data}}(\mathbf{x})}{\frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2}} \right) + \log \left(\frac{1}{2} \right) \right] + \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \left(\frac{p_\theta(\mathbf{x})}{\frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2}} \right) + \log \left(\frac{1}{2} \right) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \left[\log \left(\frac{p_{\text{data}}(\mathbf{x})}{\frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2}} \right) \right] + \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \left(\frac{p_\theta(\mathbf{x})}{\frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2}} \right) \right] + 2 \log \left(\frac{1}{2} \right) \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \left[\log \left(\frac{p_{\text{data}}(\mathbf{x})}{\frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2}} \right) \right] + \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \left(\frac{p_\theta(\mathbf{x})}{\frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2}} \right) \right] - \log 4 \\ &= D_{\text{KL}} \left(p_{\text{data}}(\mathbf{x}) \middle\| \frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2} \right) + D_{\text{KL}} \left(p_\theta(\mathbf{x}) \middle\| \frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2} \right) - \log 4 \\ &= 2 \cdot \underbrace{D_{\text{KL}} \left(p_{\text{data}}(\mathbf{x}) \middle\| \frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2} \right) + D_{\text{KL}} \left(p_\theta(\mathbf{x}) \middle\| \frac{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}{2} \right)}_{\text{Jensen-Shannon divergence (JSD)}} - \log 4 \\ &= 2D_{\text{JSD}} (p_{\text{data}}(\mathbf{x}) \| p_\theta(\mathbf{x})) - \log 4 \end{aligned}$$

Thus, the training objective becomes

$$\begin{aligned} \min_G \max_D V(G, D) &= \min_G V(G, D^*) \\ &= \min_G 2D_{\text{JSD}} (p_{\text{data}}(\mathbf{x}) \| p_\theta(\mathbf{x})) - \log 4 \end{aligned}$$

- Issue: It often suffers from vanishing gradients when the discriminator becomes too strong, making it hard for the generator to learn.

- According to the sigmoid function $\sigma(x) = \frac{1}{1+e^{-x}}$,

$$\frac{\sigma'(x)}{1 - \sigma(x)} = \frac{\frac{e^{-x}}{(1+e^{-x})^{-2}}}{1 - \frac{1}{1+e^{-x}}} = \frac{\frac{e^{-x}}{(1+e^{-x})^{-1}}}{1 + e^{-x} - 1} = \frac{\frac{e^{-x}}{(1+e^{-x})^{-1}}}{e^{-x}} = \frac{1}{(1 + e^{-x})} = \sigma(x)$$

Therefore, by calculation,

$$\begin{aligned} \frac{\partial \mathcal{L}_G^{\text{minimax}}}{\partial \theta} &= \frac{\partial}{\partial \theta} \mathbb{E}_{z \sim p(z)} [\log(1 - \sigma(h_\phi(G_\theta(z))))] \\ &= \mathbb{E}_{z \sim p(z)} \left[\frac{\partial}{\partial \theta} \log(1 - \sigma(h_\phi(G_\theta(z)))) \right] \\ &= \mathbb{E}_{z \sim p(z)} \left[-\frac{\sigma'(h_\phi(G_\theta(z)))}{1 - \sigma(h_\phi(G_\theta(z)))} \cdot \frac{\partial}{\partial \theta} h_\phi(G_\theta(z)) \right] \\ &= -\mathbb{E}_{z \sim p(z)} \left[\frac{\sigma'(h_\phi(G_\theta(z)))}{1 - \sigma(h_\phi(G_\theta(z)))} \cdot \frac{\partial}{\partial \theta} h_\phi(G_\theta(z)) \right] \\ &= -\mathbb{E}_{z \sim p(z)} \left[\sigma(h_\phi(G_\theta(z))) \cdot \frac{\partial}{\partial \theta} h_\phi(G_\theta(z)) \right] \\ &= -\mathbb{E}_{z \sim p(z)} \left[D_\phi(G_\theta(z)) \cdot \frac{\partial}{\partial \theta} h_\phi(G_\theta(z)) \right] \end{aligned}$$

When the discriminator is strong, it confidently predicts fake samples as fake, i.e., $D_\phi(G_\theta(z)) \approx 0$. Then the gradient of the generator's minimax loss is

$$\begin{aligned} \frac{\partial \mathcal{L}_G^{\text{minimax}}}{\partial \theta} &= -\mathbb{E}_{z \sim p(z)} \left[D_\phi(G_\theta(z)) \cdot \frac{\partial}{\partial \theta} h_\phi(G_\theta(z)) \right] \\ &\approx -\mathbb{E}_{z \sim p(z)} \left[0 \cdot \frac{\partial}{\partial \theta} h_\phi(G_\theta(z)) \right] \\ &= 0 \end{aligned}$$

Thus, the gradient vanishes, causing the generator to fail to improve and the training to stall, while the discriminator continues to classify generated samples as fake with high confidence.

6.12 Non-saturating GAN

- Definition:
 - Same as standard GAN
 - Generator non-saturating loss
- Note: To address the issue in the standard GAN, the non-saturating GAN uses an alternative generator loss, which encourages the generator to maximize the probability that the discriminator classifies its samples as real. Even when $D_\phi(G_\theta(\mathbf{z})) \approx 0$, the gradient remains large, providing strong learning signals and avoiding vanishing gradients. This results in stronger gradients for the generator and generally leads to more stable and faster training.

6.13 KL-based GAN

- Definition:
 - Same as standard GAN
 - Generator KL-based loss

* When the discriminator is optimal, i.e., $D_\phi(\mathbf{x}) = D_\phi^*(\mathbf{x})$,

$$\mathcal{L}_G(\theta; \phi) = D_{\text{KL}}(p_\theta(\mathbf{x}) \parallel p_{\text{data}}(\mathbf{x}))$$

meaning the generator minimizes the KL divergence between its distribution and the data distribution.

- Issue: See case 1 in Type Issues

6.14 Conditional GAN with Projection Discriminator

- Definition:
 -
 -

6.15 Wasserstein GAN (WGAN)

- Definition:
 - WGAN critic loss
 - WGAN generator loss
- Issue: See case 2 in Type Issues

6.16 Wasserstein GAN with Gradient Penalty (WGAN-GP)

- Definition:
 - WGAN-GP critic loss
 - Same as WGAN
- Note: WGAN uses the Wasserstein distance instead of JS or KL divergence. This gives smoother gradients, reduces instability, and allows more reliable convergence.

6.17 Issue Analysis

Assume that $\epsilon \rightarrow 0$ and

$$\begin{aligned} p_\theta(x) &= \mathcal{N}(x | \theta, \epsilon^2) \\ p_{\text{data}}(x) &= \mathcal{N}(x | \theta_0, \epsilon^2) \end{aligned}$$

- Case 1:

$$\begin{aligned} D_\phi : \mathbb{R} &\rightarrow [0, 1] \\ \mathcal{L}_D(\phi; \theta) &= -\mathbb{E}_{x \sim p_{\text{data}}(x)}[\log D_\phi(x)] - \mathbb{E}_{x \sim p_\theta(x)}[\log(1 - D_\phi(x))] \\ \mathcal{L}_G(\theta; \phi) &= \mathbb{E}_{x \sim p_\theta(x)}[\log(1 - D_\phi(x))] - \mathbb{E}_{x \sim p_\theta(x)}[\log(1 - D_\phi(x))] \end{aligned}$$

By calculation,

$$\begin{aligned} \mathcal{L}_G(\theta; \phi) &= D_{\text{KL}}(p_\theta(x) \| p_{\text{data}}(x)) \\ D_{\text{KL}}(p_\theta(x) \| p_{\text{data}}(x)) &= \frac{(\theta - \theta_0)^2}{2\epsilon^2} \\ \frac{\partial}{\partial \theta} D_{\text{KL}}(p_\theta(x) \| p_{\text{data}}(x)) &= \frac{\theta - \theta_0}{\epsilon^2} \end{aligned}$$

Thus, when $\epsilon \rightarrow 0$,

$$\mathcal{L}_G(\theta; \phi) \rightarrow \infty, \quad \frac{\partial}{\partial \theta} \mathcal{L}_G(\theta; \phi) \rightarrow \infty$$

This means even a tiny difference between θ and θ_0 leads to a huge generator loss and huge gradient.

When $\epsilon \rightarrow 0$, the distributions

$$\begin{aligned} p_\theta(x) &= \mathcal{N}(x | \theta, \epsilon^2) \\ p_{\text{data}}(x) &= \mathcal{N}(x | \theta_0, \epsilon^2) \end{aligned}$$

become very narrow Gaussians, each concentrated around its mean. If $\epsilon \neq \epsilon_0$, the distributions have almost no overlap.

This is problematic:

- During GAN training, the generator overreacts to small errors:
 1. The loss becomes extreme large.
 2. Even tiny change in model parameters cause big unstable updates.
- As a result, training becomes unstable and may fail to converge.

- Case 2:

$$\begin{aligned} D_\phi : \mathbb{R} &\rightarrow \mathbb{R} \\ \mathcal{L}_D(\phi; \theta) &= -\mathbb{E}_{x \sim p_{\text{data}}(x)}[D_\phi(x)] + \mathbb{E}_{x \sim p_\theta(x)}[D_\phi(x)] \\ \mathcal{L}_G(\theta; \phi) &= -\mathbb{E}_{x \sim p_\theta(x)}[D_\phi(x)] \end{aligned}$$

When $\epsilon \rightarrow 0$, both distributions collapse to single points,

$$p_\theta(x) = \theta, \quad p_{\text{data}}(x) = \theta_0$$

Then the discriminator loss becomes,

$$\mathcal{L}_D(\phi; \theta) = D_\phi(\theta) - D_\phi(\theta_0)$$

Since $D_\phi \in (-\infty, \infty)$, let

$$D_\phi(\theta) \rightarrow -\infty, \quad D_\phi(\theta_0) \rightarrow \infty$$

so that

$$\mathcal{L}_D(\phi; \theta) \rightarrow -\infty$$

This means the discriminator loss can be made arbitrarily small, i.e., there is no finite minimum and no optimal discriminator exists.

- Case 3:

$$D_\phi : \mathbb{R} \rightarrow \mathbb{R}, \quad \frac{dD_\phi}{dx} : \mathbb{R} \rightarrow [-1, 1]$$

$$\mathcal{L}_D(\phi; \theta) = -\mathbb{E}_{x \sim p_{\text{data}}(x)}[D_\phi(x)] + \mathbb{E}_{x \sim p_\theta(x)}[D_\phi(x)]$$

$$\mathcal{L}_G(\theta; \phi) = -\mathbb{E}_{x \sim p_\theta(x)}[D_\phi(x)]$$

6.18 f -divergence

f divergences

Many more f-divergences!

Name	$D_f(P\ Q)$	Generator $f(u)$
Total variation	$\frac{1}{2} \int p(x) - q(x) dx$	$\frac{1}{2} u - 1 $
Kullback-Leibler	$\int p(x) \log \frac{p(x)}{q(x)} dx$	$u \log u$
Reverse Kullback-Leibler	$\int q(x) \log \frac{q(x)}{p(x)} dx$	$-\log u$
Pearson χ^2	$\int \frac{(q(x) - p(x))^2}{p(x)} dx$	$(u - 1)^2$
Neyman χ^2	$\int \frac{(p(x) - q(x))^2}{q(x)} dx$	$\frac{(1-u)^2}{u}$
Squared Hellinger	$\int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx$	$(\sqrt{u} - 1)^2$
Jeffrey	$\int (p(x) - q(x)) \log \frac{p(x)}{q(x)} dx$	$(u - 1) \log u$
Jensen-Shannon	$\frac{1}{2} \int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx$	$-(u+1) \log \frac{1+u}{2} + u \log u$
Jensen-Shannon-weighted	$\int p(x) \pi \log \frac{\pi p(x)+(1-\pi)q(x)}{p(x)+(1-\pi)q(x)} dx + (1-\pi)q(x) \log \frac{\pi p(x)+(1-\pi)q(x)}{p(x)+(1-\pi)q(x)} dx$	$\pi u \log u - (1 - \pi + \pi u) \log(1 - \pi + \pi u)$
GAN	$\int p(x) \log \frac{p(x)+q(x)}{p(x)+q(x)} dx - \log(4)$	$u \log u - (u+1) \log(u+1)$
α -divergence ($\alpha \notin \{0, 1\}$)	$\frac{1}{\alpha(\alpha-1)} \int (p(x) [\left(\frac{q(x)}{p(x)}\right)^\alpha - 1] - \alpha(q(x) - p(x))) dx$	$\frac{1}{\alpha(\alpha-1)} (u^\alpha - 1 - \alpha(u-1))$

Source: Nowozin et al., 2016