

Unsupervised Learning & Generative AI

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

Machine Learning

Outline

1 Unsupervised Learning

- Text Models
- Image Models

2 ChatGPT

3 Autoencoders (AE)

- Manifold Learning*

4 Variational Autoencoders (VAE)

5 Flow-based Models

6 Diffusion Models

7 Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

Unsupervised Learning

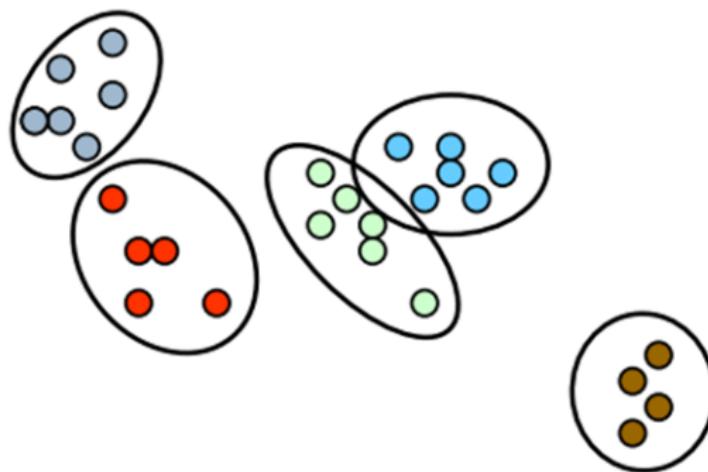
- Dataset: $\mathbb{X} = \{\mathbf{x}^{(i)}\}_i$, where $\mathbf{x}^{(i)}$'s are i.i.d. samples of \mathbf{x}
 - No supervision $\mathbf{y}^{(i)}$ (labels)
- What can we learn without labels?

Unsupervised Learning

- Dataset: $\mathbb{X} = \{\mathbf{x}^{(i)}\}_i$, where $\mathbf{x}^{(i)}$'s are i.i.d. samples of \mathbf{x}
 - No supervision $\mathbf{y}^{(i)}$ (labels)
- What can we learn without labels? The ***structures*** in \mathbb{X}
 - Inter-sample structures
 - Intra-sample structures

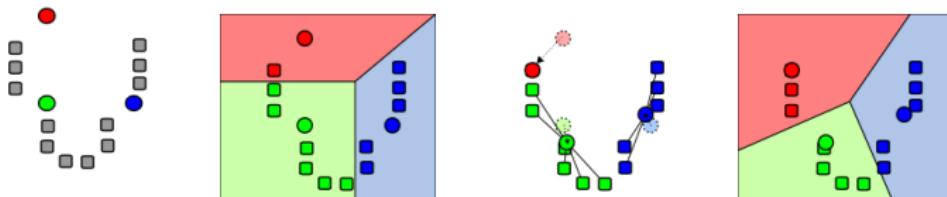
Clustering I

- Goal: to divide $x^{(i)}$'s into K groups/**clusters**
 - Based on some similarity/distance measure between $x^{(i)}$ and $x^{(j)}$



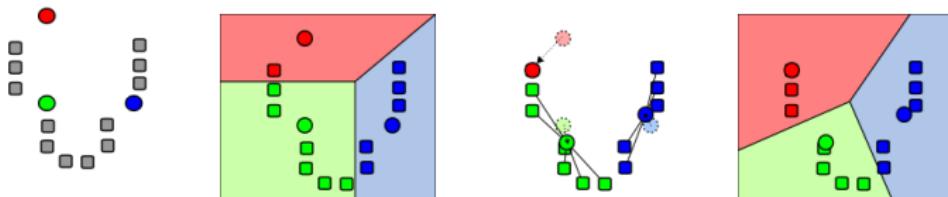
Clustering II

- K -means algorithm (K fixed): iteratively move clusterheads until convergence

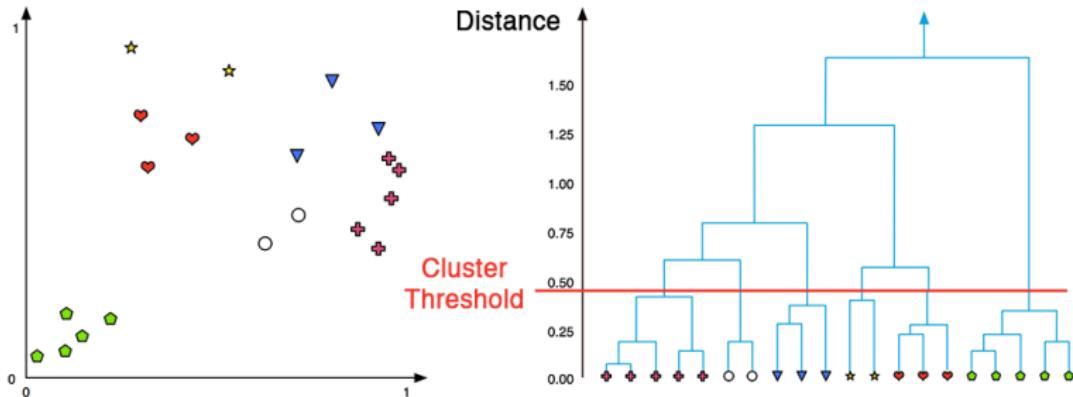


Clustering II

- K -means algorithm (K fixed): iteratively move clusterheads until convergence

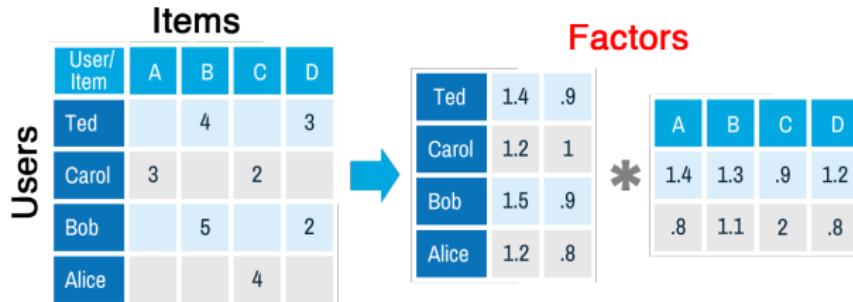


- Hierarchical clustering (variable K): iteratively merge two points/groups, then cut



Factorization for Tabular \mathbb{X}

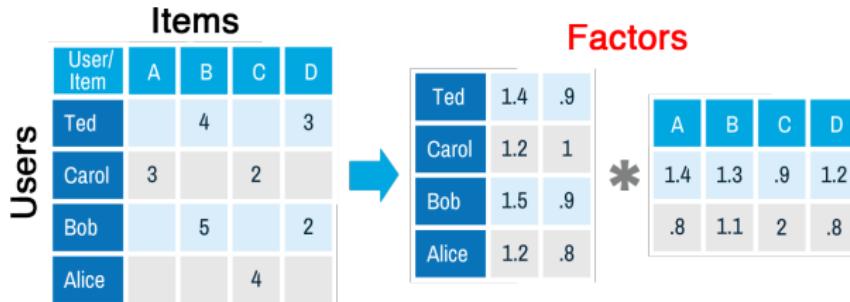
- Let X be a rating matrix where $X_{i,:} = \mathbf{x}^{(i)}$



- Goal: to approximate X with a dense matrix $\hat{X} = WH$
 - To make recommendations based on $\hat{X}_{i,:}$ for user i , known as **collaborative filtering**
- How?

Factorization for Tabular \mathbb{X}

- Let X be a rating matrix where $X_{i,:} = \mathbf{x}^{(i)}$



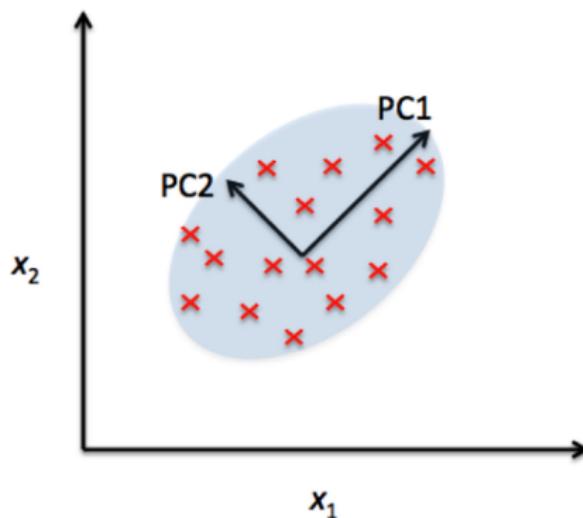
- Goal: to approximate X with a dense matrix $\hat{X} = \mathbf{WH}$
 - To make recommendations based on $\hat{X}_{i,:}$ for user i , known as **collaborative filtering**
- How? Non-negative matrix factorization (NMF) [19, 20]:

$$\arg \min_{\mathbf{W} \geq \mathbf{0}, \mathbf{H} \geq \mathbf{0}} \|\mathbf{X} - \mathbf{WH}\|_F$$

- So, positive elements in \mathbf{W} and \mathbf{H} can be seen as **factor** degrees

Dimension Reduction

- Goal: to learn a low dimensional representation \mathbf{z} of \mathbf{x}
 - E.g., PCA



Self-Supervised Learning

- Goal: to learn a model that is able to “fill in the blanks”



Self-Supervised Learning

- Goal: to learn a model that is able to “fill in the blanks”
- Links unsupervised tasks with supervised models
 - Much more training data for models
 - Representation learning with deep networks

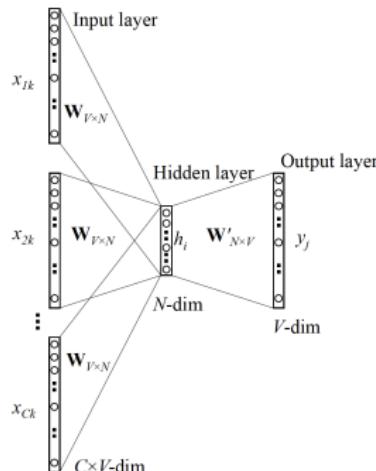


Example I: Word2Vec

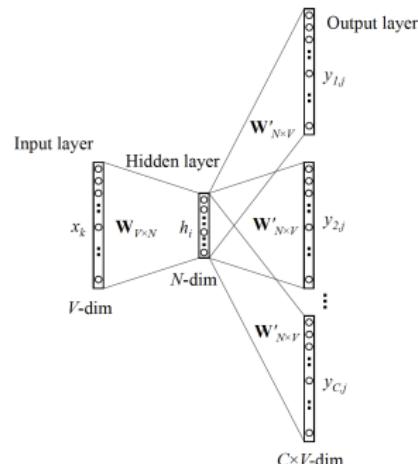
- Goal: to learn a model for blank filling

Example I: Word2Vec

- Goal: to learn a model for blank filling
 - E.g., word2vec [25, 24]: "... *the cat sat on...*"



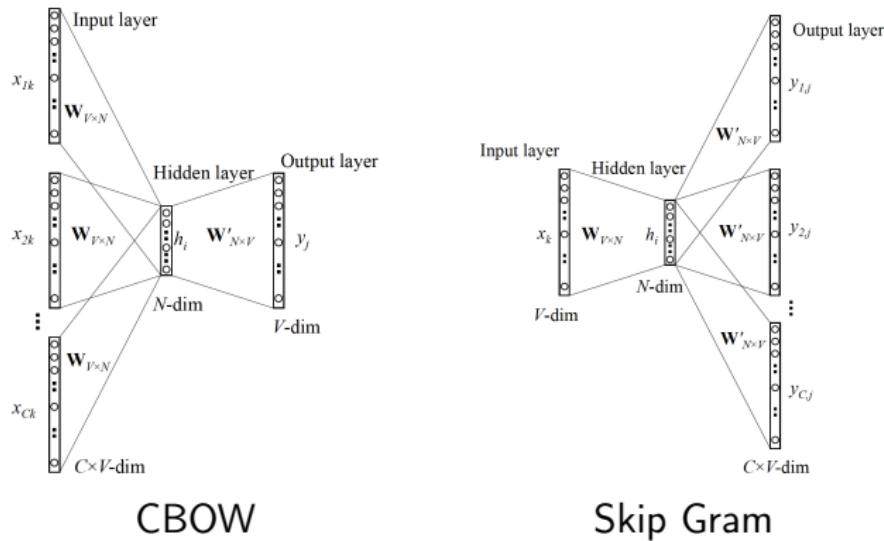
CBOW



Skip Gram

Example I: Word2Vec

- Goal: to learn a model for blank filling
 - E.g., word2vec [25, 24]: "... the cat sat on..."



- Latent representation \mathbf{h} encodes the ***semantics*** of a word
 - No need for synonym dictionary; big data tell that already

Example II: Doc2Vec

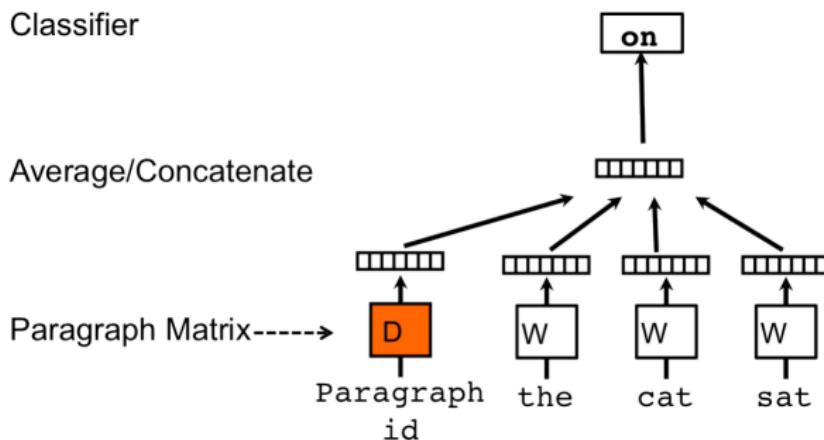
- How to encode a document?

Example II: Doc2Vec

- How to encode a document?
- Bag of words (TF-IDF), average word2vec, etc.
 - Do **not** capture the semantics due to sentence/paragraph/doc structure
 - "*John likes Mary*" \neq "*Mary likes John*"

Example II: Doc2Vec

- How to encode a document?
- Bag of words (TF-IDF), average word2vec, etc.
 - Do **not** capture the semantics due to sentence/paragraph/doc structure
 - “John likes Mary” ≠ “Mary likes John”
- Why not apply self-supervised learning to docs?
 - Doc2vec [17]: to capture the **context** not explained by words
 - **Transductive** rather than inductive; does not work with unseen docs



Generative Models

- Goal: to generate new samples of x
 - Can be conditioned on instructions (input)
- Largely based on self-supervised learning

The image is a collage of four screenshots of AI interfaces:

- DALL-E 2:** A dark-themed interface with the text "DALL-E 2" and "DALL-E is an AI system that can create realistic images and art from descriptions or natural language".
- ChatGPT:** A white-themed interface with the text "Welcome to ChatGPT" and "Log in with your OpenAI account to continue". It features a large "G" logo.
- Adobe Firefly:** A white-themed interface with the text "Text to image" and "Adobe Firefly". It shows examples of generated images like a lighthouse and a colorful abstract scene.
- Bard:** A white-themed interface with the text "Bard can explain why large language models might make mistakes". It includes a "Meet Bard" section and a note that "Bard isn't currently supported in your country. Try again". It features a large "G" logo.

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

Generating Text

- With large (self-supervised) training data, *transformers* [40] have shown to perform better than RNNs and CNNs
 - Mainly due to the $O(1)$ point distance
- Two common text models based on transformer:
 - BERT [7]: non-autoregressive
 - GPT [30]: autoregressive

BERT [7]

- Massively trained *encoder* of the original transformer
 - Non-autoregressive
- Pre-training tasks:
 - Masked language model (“Cloze” task)
 - “A quick brown [MASK] jumps over the lazy dog” → “fox” 11%, “ant” 5%, ...
 - Next sentence prediction
 - “[MASK] go to store [SEP] to buy a [MASK] of milk” → True 93%, False 7%

One Pre-training, Multiple Fine-tuning Tasks

- Special input to identify downstream task: [CLS] token

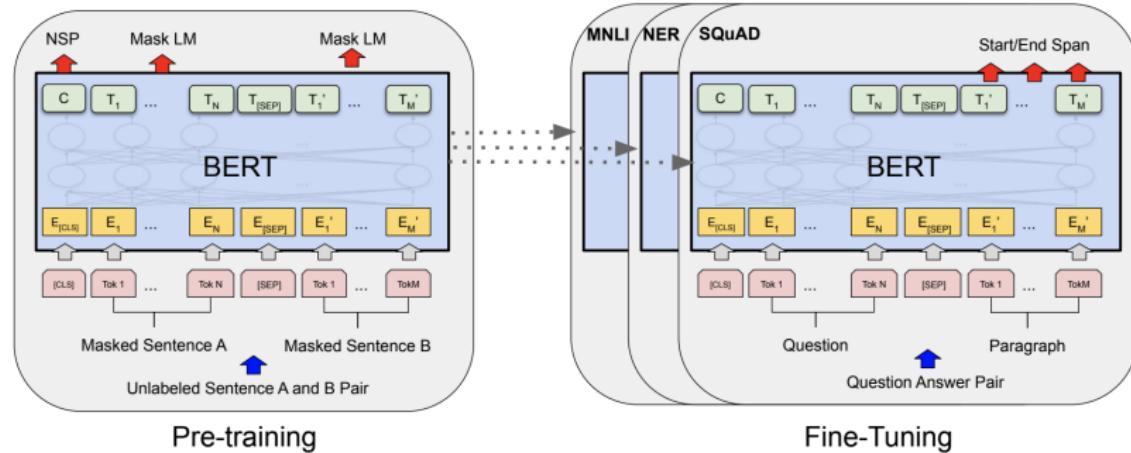
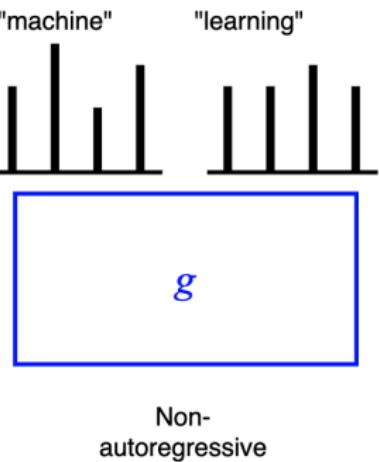
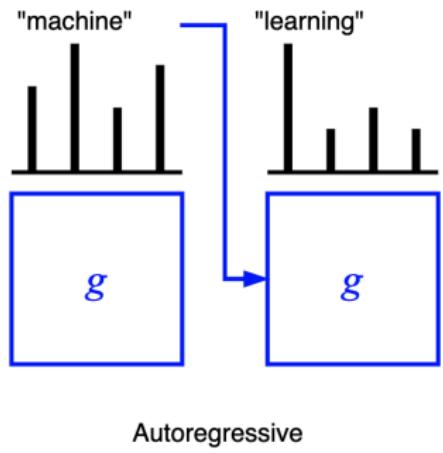


Figure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the same architectures are used in both pre-training and fine-tuning. The same pre-trained model parameters are used to initialize models for different down-stream tasks. During fine-tuning, all parameters are fine-tuned. [CLS] is a special symbol added in front of every input example, and [SEP] is a special separator token (e.g. separating questions/answers).

GPT [30]

- Massively trained **decoder** of the original transformer
 - Autoregressive
- Multitask pre-training by maximizing $\text{Pr}(\text{output}|\text{input}, \text{task})$:
 - Translation: $\text{Pr}(\text{french text}|\text{en text}, \text{translation})$
 - Question answering: $\text{Pr}(\text{answer}|\text{question}, \text{qa})$
 - Reading comprehension: $\text{Pr}(\text{answer}|\text{document}, \text{question}, \text{reading})$
- Usage for downstream task: fine-tuning or ***prompting***

Autoregressive or Not?



- It's easier for an autoregressive model to generate coherent text
- In this lecture, we focus on GPT and its variants

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

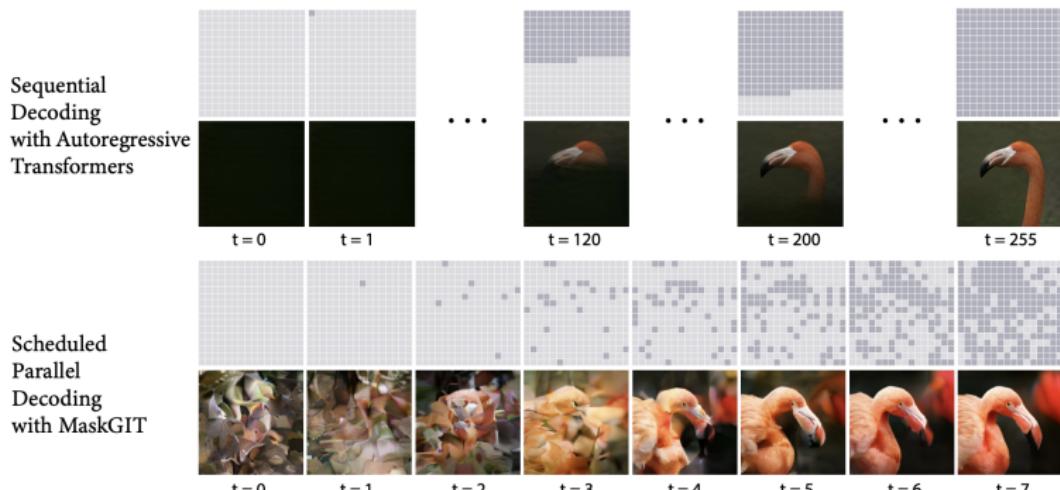
⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

Autoregressive or Not?

- Autoregressive models can still give good (if not better) performance
 - E.g., PixelRNN [39], PixelCNN [38], or MaskGIT [5]
- But unlike in text domain, generating images pixel-by-pixel is **very slow**
- Speed-up?
 - Scheduled **parallel-pixel** generation (e.g., MaskGIT)
 - **Stepwise** generation of **whole** images (e.g., flow-based and diffusion models)



Whole-Image Generation I

- Goal: given \mathbb{X} , to learn a *generator function* g such that $\hat{\mathbf{x}} = g(\mathbf{c}; \Theta_g)$ looks like a real image in \mathbb{X}
 - \mathbf{c} is a code or condition
 - Θ_g represents parameters of g
- Objective from Information Theory perspective:

$$\arg \min_g D(P_{\text{data}} \| P_g)$$

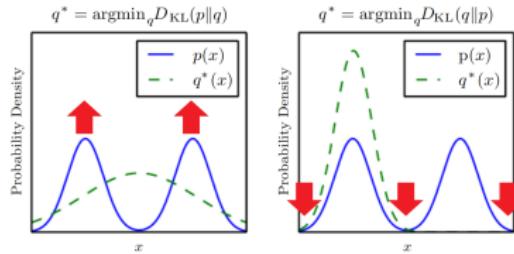
- P_{data} is the distribution of real data in the ground truth
- P_g is the distribution of generated data
- D is a divergence measure, e.g., D_{KL}

Whole-Image Generation I

- Goal: given \mathbb{X} , to learn a **generator function g** such that $\hat{\mathbf{x}} = g(\mathbf{c}; \Theta_g)$ looks like a real image in \mathbb{X}
 - \mathbf{c} is a code or condition
 - Θ_g represents parameters of g
- Objective from Information Theory perspective:

$$\arg \min_g D(P_{\text{data}} \| P_g)$$

- P_{data} is the distribution of real data in the ground truth
 - P_g is the distribution of generated data
 - D is a divergence measure, e.g., D_{KL}
- Why not $D_{\text{KL}}(P_g \| P_{\text{data}})$?



Whole-Image Generation II

- Minimizing $D_{KL}(P_{\text{data}} \| P_g)$ amounts to maximizing $P(\mathbb{X} | \Theta_g)$, the log likelihood of Θ_g :

$$\begin{aligned} g^* &= \arg \min_g D_{KL}(P_{\text{data}} \| P_g) \\ &= \arg \max_g E_{\mathbf{x} \sim P_{\text{data}}} [\log P_g(\mathbf{x})] + H(\mathbf{x} \sim P_{\text{data}}) \\ &= \arg \max_g E_{\mathbf{x} \sim P_{\text{data}}} [\log P_g(\mathbf{x})] \\ &\approx \arg \max_g \sum_{\mathbf{x}^{(i)} \in \mathbb{X}} \log P_g(\mathbf{x}^{(i)}) \\ &= \arg \max_g \log \prod_{\mathbf{x}^{(i)} \in \mathbb{X}} P_g(\mathbf{x}^{(i)}) \end{aligned}$$

$$\begin{aligned} \Theta_g^* &= \arg \max_{\Theta_g} \log \prod_{\mathbf{x}^{(i)} \in \mathbb{X}} P(\mathbf{x}^{(i)} | \Theta_g) \\ &= \arg \max_{\Theta_g} \log P(\mathbb{X} | \Theta_g) \end{aligned}$$

- Other divergence measures can lead to similar results

Common Whole-Image Generation Methods

- Autoencoder (single-step: $\hat{\mathbf{x}} = g(\mathbf{c})$)
 - Pros: easy
 - Cons: but no creativity, blurry images
- Variational Autoencoder (single-step: $\hat{\mathbf{x}} = g(\mathbf{c})$)
 - Pros: creative
 - Cons: only maximizes a lower bound of $P(\mathbb{X}|\Theta_g)$; blurry images
- Flow-based methods (multi-step: $\hat{\mathbf{x}} = g^{(T)}(\cdots g^{(2)}(g^{(1)}(\mathbf{c})))$)
 - Pros: maximizes $P(\mathbb{X}|\Theta_g)$ directly
 - Cons: limited expressiveness of g for invertibility; slow training and inference
- GANs (single-step: $\hat{\mathbf{x}} = g(\mathbf{c})$)
 - Pros: good image quality (sharp and coherent)
 - Cons: difficult to train (convergence issue, mode collapse, vanishing gradients, etc.)
- Diffusion models (multi-step: $\hat{\mathbf{x}} = g^{(T)}(\cdots g^{(2)}(g^{(1)}(\mathbf{c})))$)
 - Pros: good image quality; efficient & stable to train, can be made conditional easily
 - Cons: slow inference

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

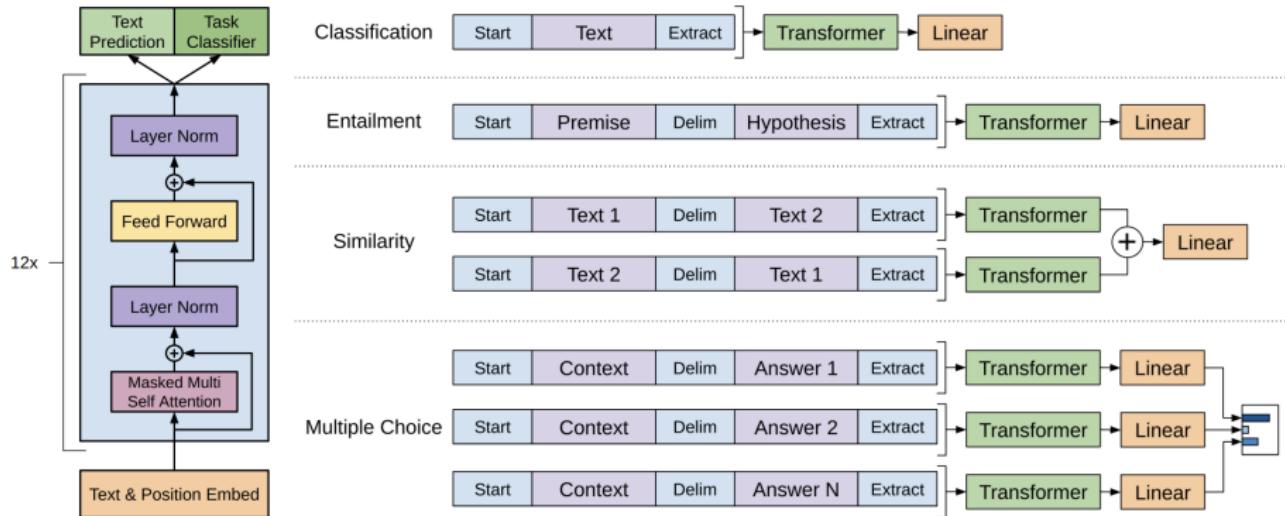
Evolutions

- 2018 GPTv1 [30]
 - Self-supervised pre-training
- 2019 GPTv2 [31]
 - **Multitask** pre-training
- 2020 GPTv3 [4]
 - **Few-shot** & **in-context** learning
- 2022 GPTv3.5 [28]
 - **Alignment** using (supervised) instruction tuning + reinforcement learning from human feedback (RLHF)
- 2023 GPT4: mixture of experts

GPTv1 [30]

- Aims at 2-step training process:
 - ① Self-supervised pre-training on unlabeled data
 - To predict next word in a sentence (language model)
 - ② Discriminative fine-tuning on labeled data in downstream tasks
 - Here, fine-tuning could also mean training a new model based on extracted features

Example Fine-tuning Tasks



Does Self-supervised Pre-training Help?

Table 5: Analysis of various model ablations on different tasks. Avg. score is a unweighted average of all the results. (*mc*= Mathews correlation, *acc*=Accuracy, *pc*=Pearson correlation)

Method	Avg. Score	CoLA (mc)	SST2 (acc)	MRPC (F1)	STS-B (pc)	QQP (F1)	MNLI (acc)	QNLI (acc)	RTE (acc)
Transformer w/ aux LM (full)	74.7	45.4	91.3	82.3	82.0	70.3	81.8	88.1	56.0
Transformer w/o pre-training	59.9	18.9	84.0	79.4	30.9	65.5	75.7	71.2	53.8
Transformer w/o aux LM	75.0	47.9	92.0	84.9	83.2	69.8	81.1	86.9	54.4
LSTM w/ aux LM	69.1	30.3	90.5	83.2	71.8	68.1	73.7	81.1	54.6

GPTv2 [31]

- Many NLP tasks can be formulated as the problem of maximizing $\text{Pr}(\text{output}|\text{input}, \text{task})$
 - Translation: $\text{Pr}(\text{french text}|\text{en text}, \text{translation})$
 - Question answering: $\text{Pr}(\text{answer}|\text{question}, \text{qa})$
 - Reading comprehension: $\text{Pr}(\text{answer}|\text{document}, \text{question}, \text{reading})$
- GPTv3 [4]: ***few-shot & in-context*** learning
- GPTv3.5 [28]: ***alignment*** using (supervised) instruction tuning & reinforcement learning from human feedback (RLHF)
- GPT4: mixture of experts

GPTv3 [4]

- The model learns from few shots (examples), even *in context*
 - Enables **prompting** techniques for downstream tasks
 - E.g., few shots, chain of thought (CoT), or simply “Let’s work this out step-by-step to ensure the answer is correct”

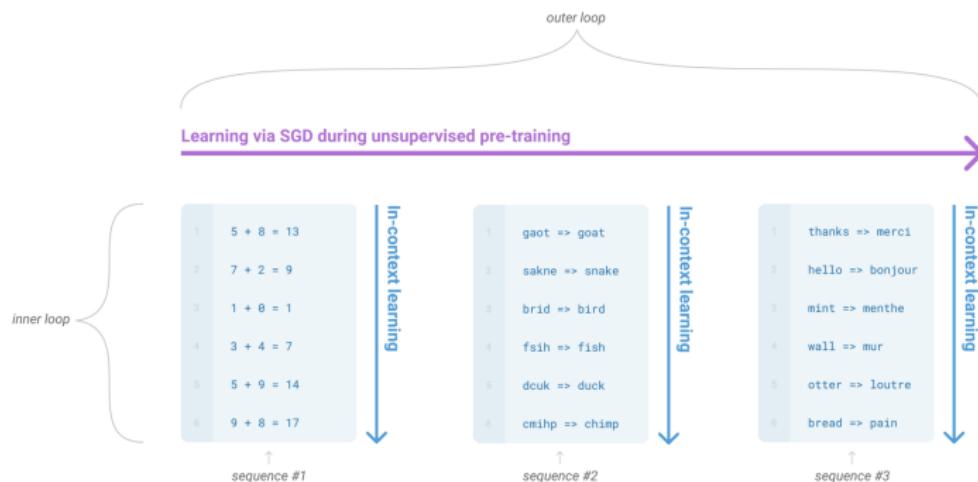


Figure 1.1: Language model meta-learning. During unsupervised pre-training, a language model develops a broad set of skills and pattern recognition abilities. It then uses these abilities at inference time to rapidly adapt to or recognize the desired task. We use the term “in-context learning” to describe the inner loop of this process, which occurs within the forward-pass upon each sequence. The sequences in this diagram are not intended to be representative of the data a model would see during pre-training, but are intended to show that there are sometimes repeated sub-tasks embedded within a single sequence.

Zero/Few Shot Prompting vs. Fine-tuning

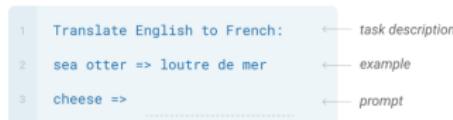
Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



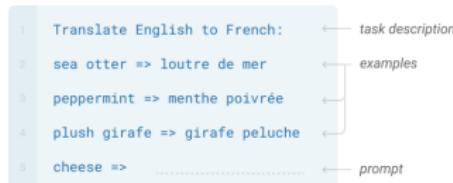
One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



Performance vs. Model Size

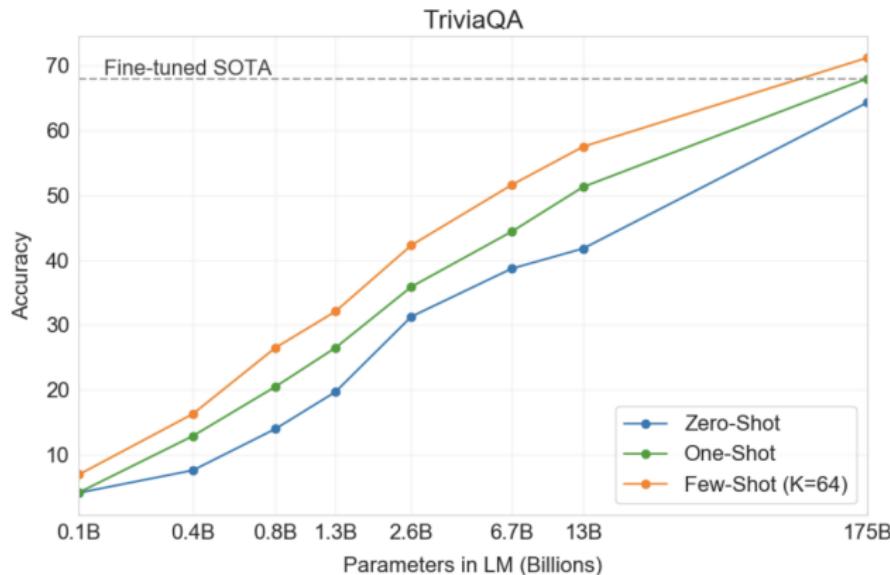


Figure 3.3: On TriviaQA GPT3's performance grows smoothly with model size, suggesting that language models continue to absorb knowledge as their capacity increases. One-shot and few-shot performance make significant gains over zero-shot behavior, matching and exceeding the performance of the SOTA fine-tuned open-domain model, RAG [LPP⁺20]

GPTv3.5 [28]

- From GPT to “ChatGPT” through *alignment*
 - Supervised (multi-task) instruction tuning
 - Reinforcement learning from human feedback (RLHF)

Step 1

Collect demonstration data, and train a supervised policy.

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old



Some people went to the moon...

A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3 with supervised learning.



Step 2

Collect comparison data, and train a reward model.

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A Explain gravity...
B Explain air...
C Moon is natural satellite of...
D People went to the moon...



D > C > A = B

A labeler ranks the outputs from best to worst.

This data is used to train our reward model.

RM
D > C > A = B

Step 3

Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.

Write a story about frogs



PPO

The policy generates an output.

Once upon a time...

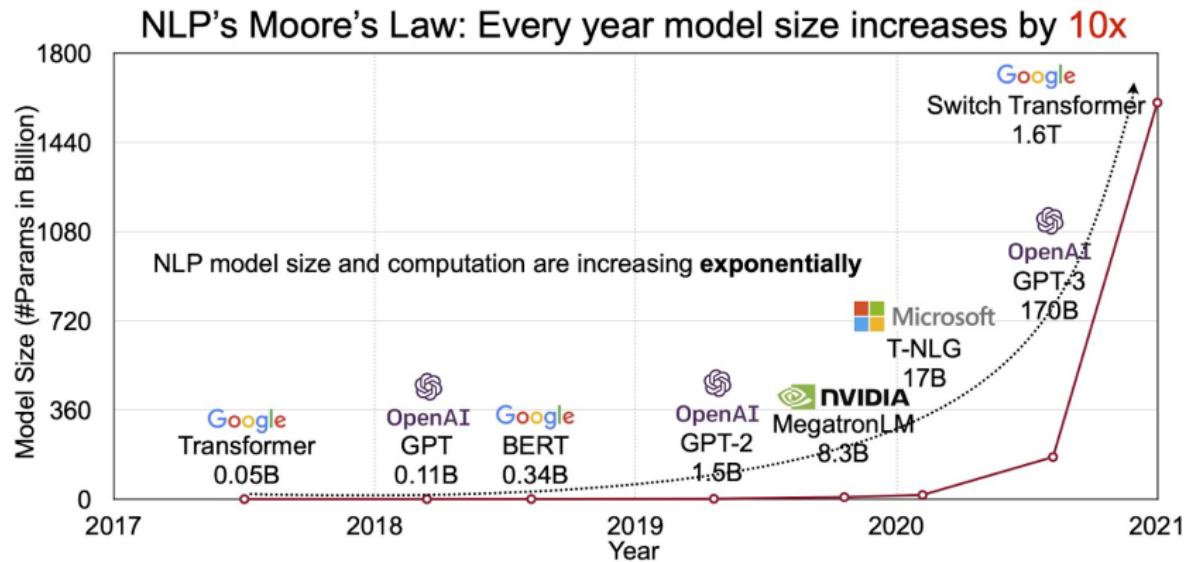


The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

r_k

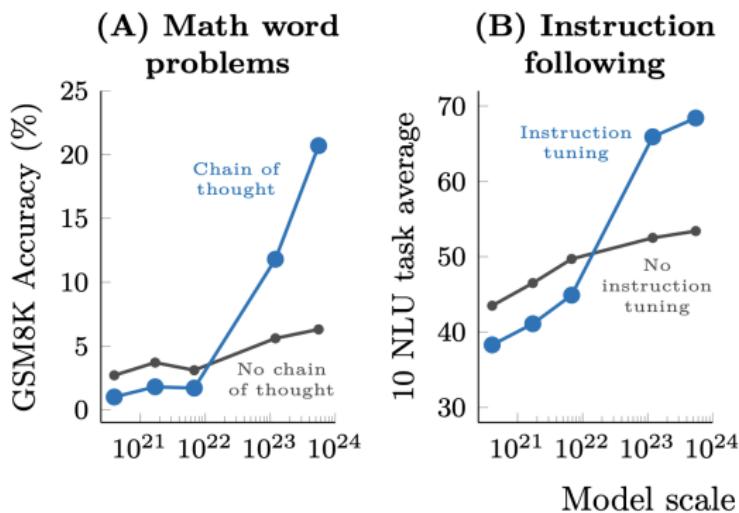
Sizes of Large Language Models (LLMs)



- Training costs [36]:
 - 110M params: \$2.5k–\$50k
 - 340M params: \$10k–\$200k
 - 1.5B param: \$80k–\$1.6m

Size Does Matter!

- Emerging abilities of LLMs [42]



- A balance: 70B parameters + 1.4T training tokens [11]

GPT Variants

- Low-rank adaptation (LoRA) [12]: to efficiently fine-tune LLMs
- WebGPT [27]: GPT that can search the web
- Retrieval-augmented generation (RAG) [21]: GPT to query external knowledge (vector) base
- GPTs that can reflect on their answers [2, 26]

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

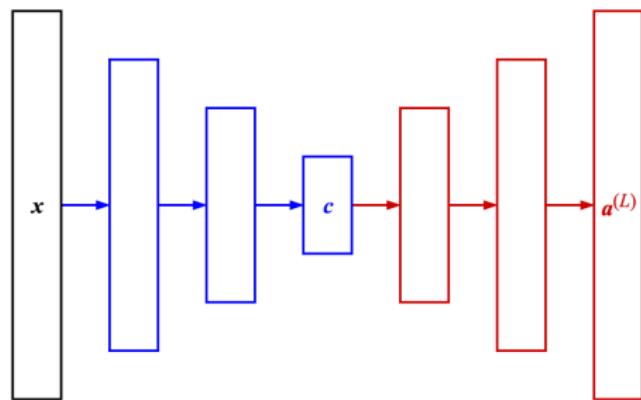
⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

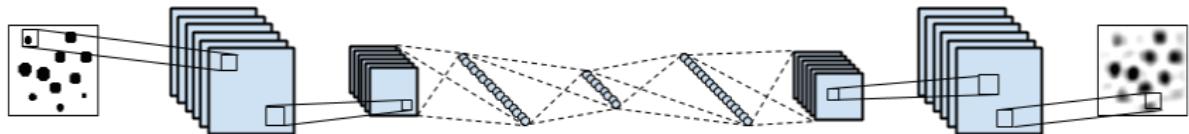
Autoencoders (AE)

- **Encoder**: to learn a low dimensional representation \mathbf{c} (called **code**) of input \mathbf{x}
- **Decoder**: to reconstruct \mathbf{x} from \mathbf{c}
- Objective: $\arg \max_{\Theta} \log P(\mathbb{X} | \Theta) = \arg \max_{\Theta} \sum_i \log P(\mathbf{x}^{(i)} | \Theta)$
- Assuming that $\mathbf{x} \sim \mathcal{N}(\mu, \cdot)$, we have linear output units
 $\mathbf{a}^{(L)} = \mathbf{z}^{(L)} = \hat{\mu}$
 - $\log P(\mathbf{x}^{(i)} | \Theta) \propto -\|\mathbf{x}^{(i)} - \mathbf{a}^{(i,L)}\|^2$
 - $\arg \max_{\Theta} \sum_i \log P(\mathbf{x}^{(i)} | \Theta) = \arg \min_{\Theta} \sum_i \|\mathbf{x}^{(i)} - \mathbf{a}^{(i,L)}\|^2$ (minimizing reconstruct error)



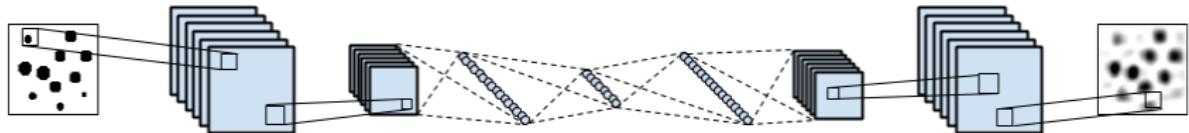
Convolutional Autoencoders

- Convolution + deconvolution layers:



Convolutional Autoencoders

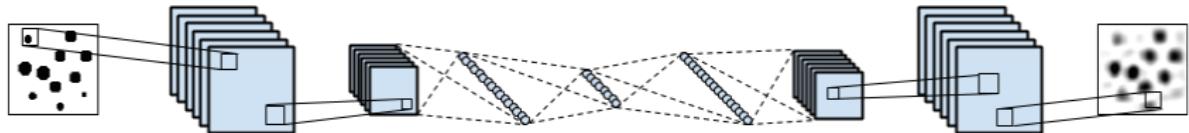
- Convolution + deconvolution layers:



- Decoder is a simplified DeconvNet [43] trained from scratch:

Convolutional Autoencoders

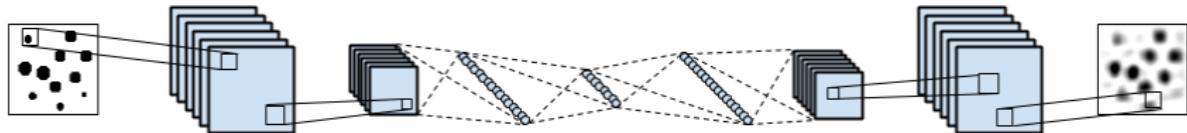
- Convolution + deconvolution layers:



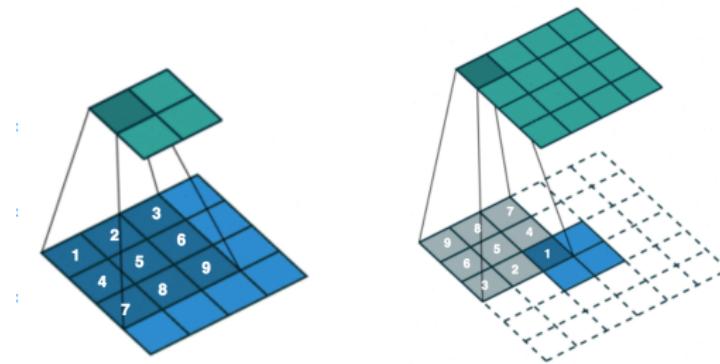
- Decoder is a simplified DeconvNet [43] trained from scratch:
 - Uppooling → upsampling (no need to remember max positions)

Convolutional Autoencoders

- Convolution + deconvolution layers:



- Decoder is a simplified DeconvNet [43] trained from scratch:
 - Uppooling → upsampling (no need to remember max positions)
 - Deconvolution → convolution



Codes & Reconstructed x

- A 32-bit code can roughly represents a 32×32 (1024 dimensional) MNIST image



Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

⑦ Generative Adversarial Networks*

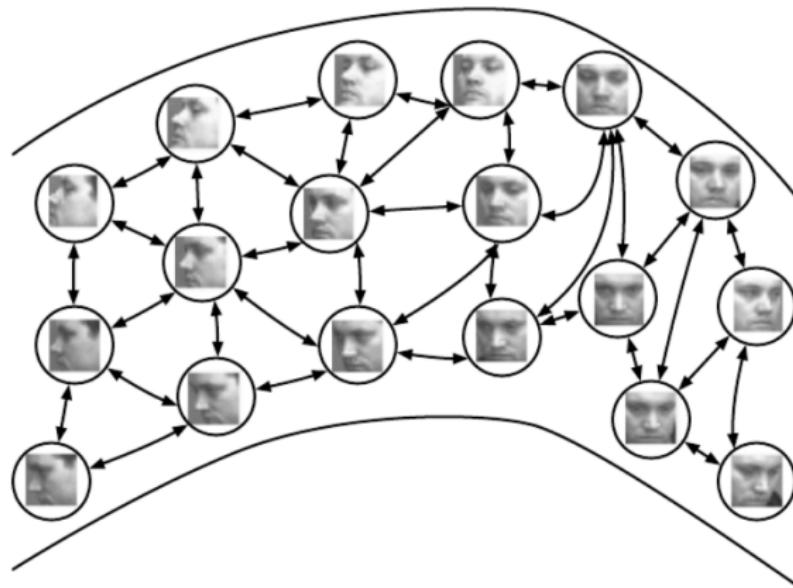
- Basic Architecture
- Challenges
- More GANs

Manifolds I

- In many applications, data concentrate around one or more low-dimensional *manifolds*

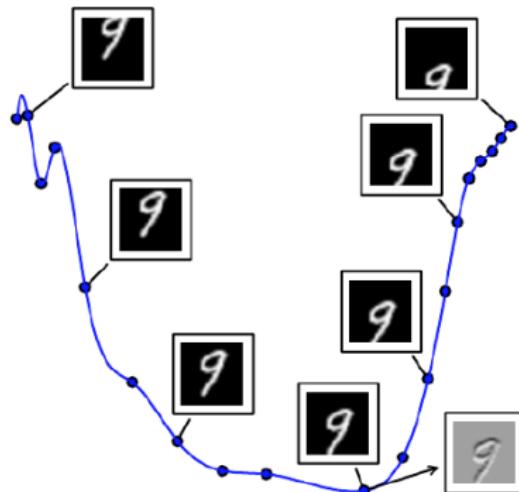
Manifolds I

- In many applications, data concentrate around one or more low-dimensional *manifolds*
- A manifold is a topological space that are *linear locally*



Manifolds II

- For each point x on a manifold, we have its **tangent space** spanned by **tangent vectors**
 - Local directions specify how one can change x infinitesimally while staying on the manifold



Learning Manifolds I

- How to make \mathbf{c} produced by autoencoders denote a *coordinate* of a dimensional manifold?

Learning Manifolds I

- How to make \mathbf{c} produced by autoencoders denote a *coordinate* of a dimensional manifold?
- Contractive autoencoder [33]: regularizes the code \mathbf{c} such that it is invariant to local changes of \mathbf{x} :

$$\Omega(\mathbf{c}) = \sum_i \left\| \frac{\partial \mathbf{c}^{(i)}}{\partial \mathbf{x}^{(i)}} \right\|_F^2$$

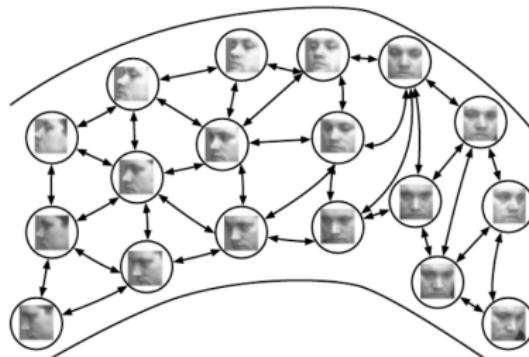
- $\partial \mathbf{c}^{(i)} / \partial \mathbf{x}^{(i)}$ is a Jacobian matrix

Learning Manifolds I

- How to make \mathbf{c} produced by autoencoders denote a *coordinate* of a dimensional manifold?
- Contractive autoencoder [33]: regularizes the code \mathbf{c} such that it is invariant to local changes of \mathbf{x} :

$$\Omega(\mathbf{c}) = \sum_i \left\| \frac{\partial \mathbf{c}^{(i)}}{\partial \mathbf{x}^{(i)}} \right\|_F^2$$

- $\partial \mathbf{c}^{(i)} / \partial \mathbf{x}^{(i)}$ is a Jacobian matrix
- Hence, \mathbf{c} represents only the variations needed to reconstruct \mathbf{x}
 - I.e., \mathbf{c} changes most along tangent vectors



Learning Manifolds II

- In practice, it is easier to train a denoising autoencoder [41]:



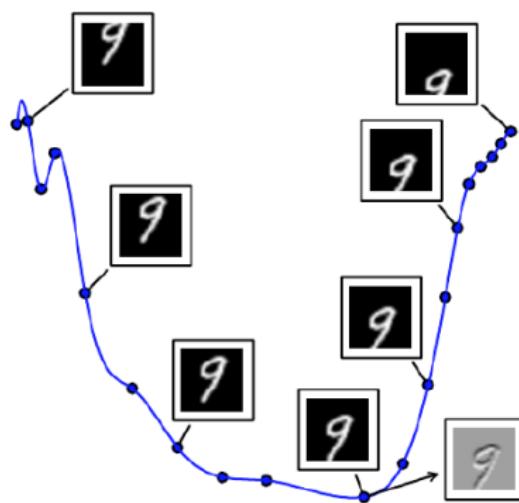
Learning Manifolds II

- In practice, it is easier to train a denoising autoencoder [41]:
 - Encoder: to encode \mathbf{x} **with** random noises
 - Decoder: to reconstruct \mathbf{x} **without** noises



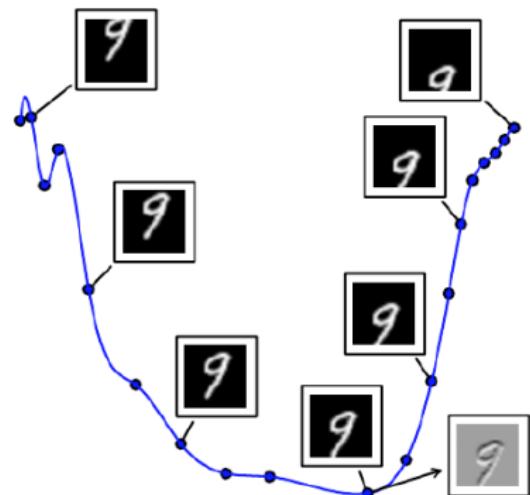
Getting Tangent Vectors I

- The code c represents a coordinate on a low dimensional manifold
 - E.g., the blue line
- How to get the tangent vectors of a given c ?



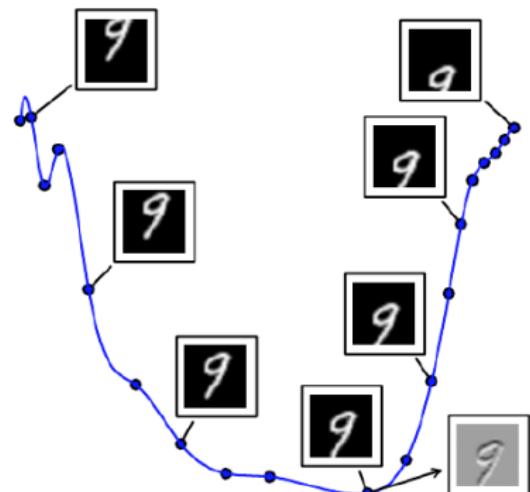
Getting Tangent Vectors II

- Recall: directions in the input space that *changes c most* should be tangent vectors



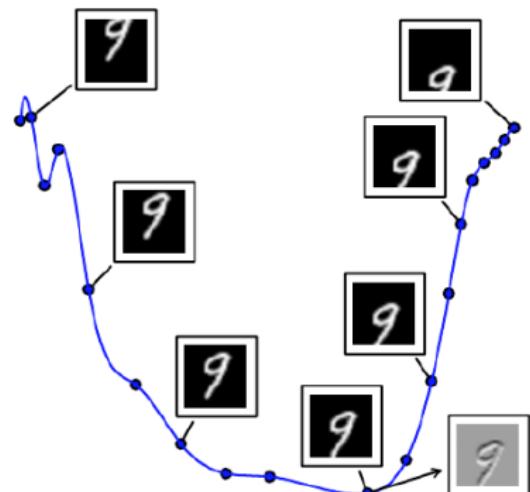
Getting Tangent Vectors II

- Recall: directions in the input space that *changes c most* should be tangent vectors
- Given a point x , let c be the code of x and $J(x) = \frac{\partial c}{\partial x}$ be the Jacobian matrix of c at x



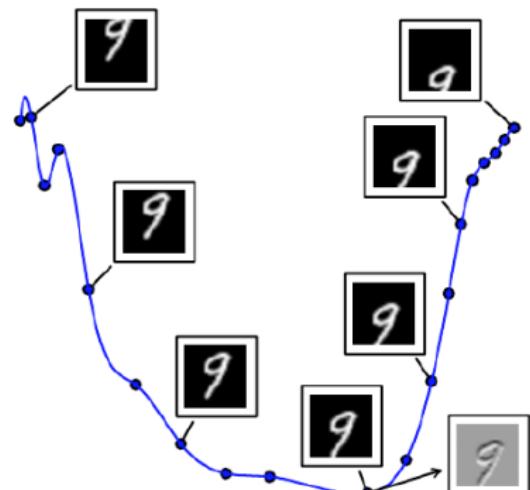
Getting Tangent Vectors II

- Recall: directions in the input space that *changes c most* should be tangent vectors
- Given a point x , let c be the code of x and $J(x) = \frac{\partial c}{\partial x}$ be the Jacobian matrix of c at x
 - $J(x)$ summarizes how c changes in terms of x



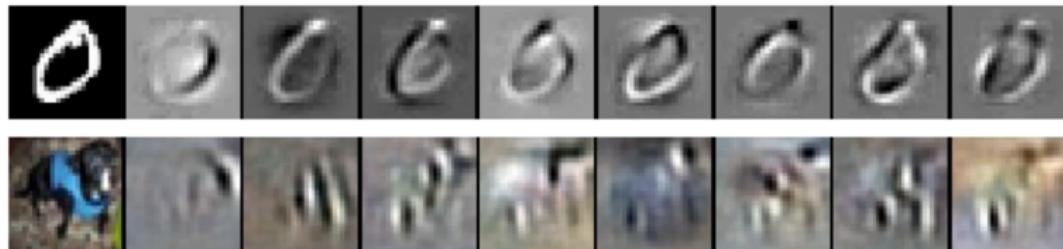
Getting Tangent Vectors II

- Recall: directions in the input space that *changes c most* should be tangent vectors
 - Given a point x , let c be the code of x and $J(x) = \frac{\partial c}{\partial x}$ be the Jacobian matrix of c at x
 - $J(x)$ summarizes how c changes in terms of x
- Decompose $J(x)$ using SVD such that $J(x) = UDV^\top$
 - Let tangent vectors be *rows of V corresponding to the largest singular values in D*



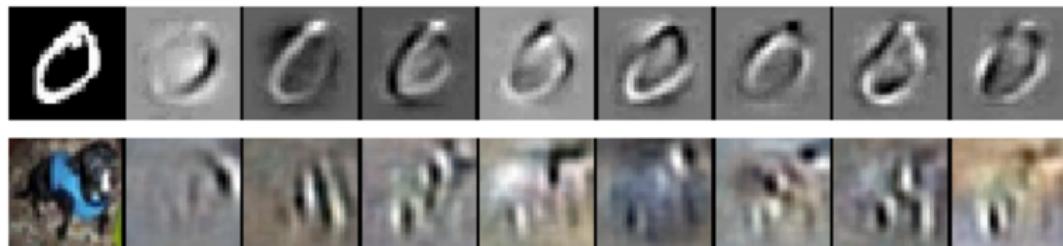
Getting Tangent Vectors III

- In practice, $J(x)$ usually has few large singular values
- Tangent vectors found by contractive/denoising autoencoders:



Getting Tangent Vectors III

- In practice, $J(\mathbf{x})$ usually has few large singular values
- Tangent vectors found by contractive/denoising autoencoders:



- Can be used by Tangent Prop [37]:
- Let $\{\mathbf{v}^{(i,j)}\}_j$ be tangent vectors of each example $\mathbf{x}^{(i)}$
- Trains an NN classifier f with cost penalty: $\Omega[f] = \sum_{i,j} \nabla_{\mathbf{x}} f(\mathbf{x}^{(i)})^\top \mathbf{v}^{(i,j)}$
 - Points in the same manifold share the same label

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

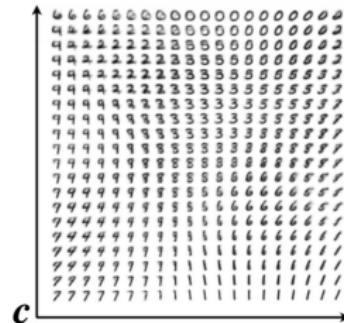
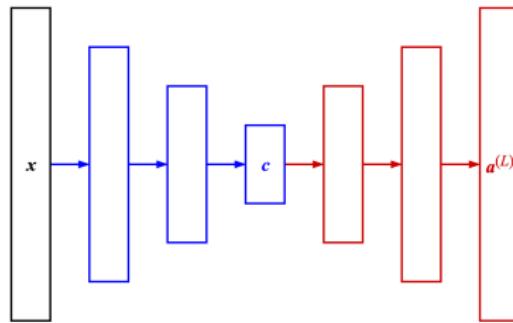
⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

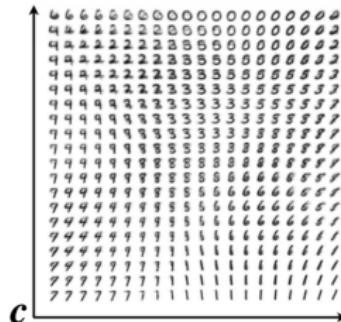
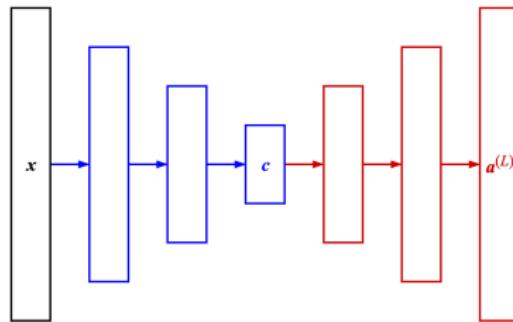
Problems of Autoencoder I

- Ideally, the decoder of an autoencoder can be used to generate images even with *synthetic codes*



Problems of Autoencoder I

- Ideally, the decoder of an autoencoder can be used to generate images even with *synthetic codes*



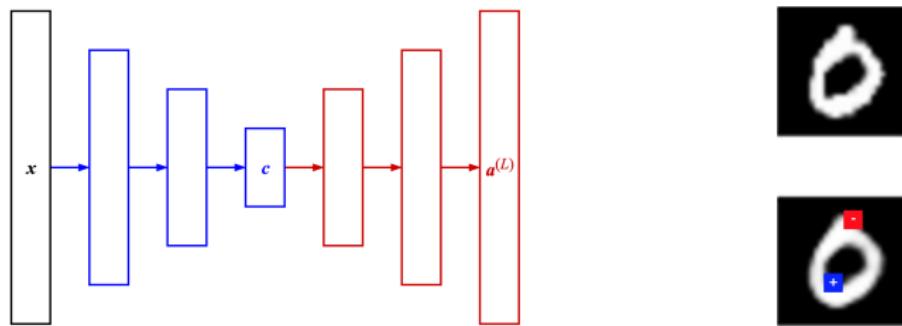
- In reality, the learnt c in code space has many “holes” that fail to map to images
 - More complex image patterns
 - Training data are never enough

Problems of Autoencoder II

- Blurry images: the objective

$$\arg \max_{\Theta} \sum_i \log P(\mathbf{x}^{(i)} | \Theta) = \arg \min_{\Theta} \sum_i \|\mathbf{x}^{(i)} - \mathbf{a}^{(i,L)}\|^2$$

does not penalize Gaussian pixel noises in $\mathbf{a}^{(i,L)}$



Variational Autoencoders (VAE) [15]

- Encoder $f(\cdot; \Theta_f)$: maps each sample of \mathbf{x} (i.e., $\mathbf{x}^{(i)} \in \mathbb{X}$) to an ***axis-aligned normal distribution*** $\mathcal{N}(\mu, \sigma)$
 - Each code dimension is independent with each other
 - $f(\mathbf{x}) = (\mu, \sigma)$
- Decoder $g(\cdot; \Theta_g)$: same as that of AE
- How to minimize the objective $\arg \max_{\Theta_f, \Theta_g} \log P(\mathbb{X} | \Theta_f, \Theta_g))$?

VAE Objective

- Objective:

$$\arg \max_{\Theta_f, \Theta_g} \log P(\mathbb{X} | \Theta_f, \Theta_g) = \arg \max_{\Theta_f, \Theta_g} \sum_i \log P(x^{(i)} | \Theta_f, \Theta_g)$$

- Considering $\log P(\mathbf{x})$ for any sample \mathbf{x} , we have

$$\begin{aligned}\log P(\mathbf{x}) &= \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log P(\mathbf{x}) d\mathbf{c} \quad // \text{Q can be any distribution} \\ &= \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{c}, \mathbf{x})}{P(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} = \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{c}, \mathbf{x})}{Q(\mathbf{c}|\mathbf{x})} \frac{Q(\mathbf{c}|\mathbf{x})}{P(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} \\ &= \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{c}, \mathbf{x})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} + \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{Q(\mathbf{c}|\mathbf{x})}{P(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} \\ &= \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{c}, \mathbf{x})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} + D_{\text{KL}}(Q(\mathbf{c}|\mathbf{x}) \| P(\mathbf{c}|\mathbf{x})) \\ &\geq \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{c}, \mathbf{x})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} \quad // \text{lower bound} \\ &= \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{x}|\mathbf{c})P(\mathbf{c})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c}\end{aligned}$$

- VAE lets $Q(\cdot|\mathbf{x}) = \mathcal{N}(f(\mathbf{x}; \Theta_f))$ and $P(\cdot|\mathbf{c}) = \mathcal{N}(g(\mathbf{c}; \Theta_g))$
 - So $Q(\mathbf{c}|\mathbf{x}) = Q(\mathbf{c}|\mathbf{x}, \Theta_f)$ and $P(\mathbf{x}|\mathbf{c}) = P(\mathbf{x}|\mathbf{c}, \Theta_g)$
- New objective: finds Θ_f and Θ_g that maximize the lower bound
 - Not necessarily maximize $\log P(\mathbf{x}|\Theta_f, \Theta_g)$

Maximizing Lower Bound

$$\begin{aligned} & \int_c Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{x}|\mathbf{c})P(\mathbf{c})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} \\ &= \int_c Q(\mathbf{c}|\mathbf{x}, \Theta_f) \log \left(\frac{P(\mathbf{x}|\mathbf{c}, \Theta_g)P(\mathbf{c})}{Q(\mathbf{c}|\mathbf{x}, \Theta_f)} \right) d\mathbf{c} \\ &= \int_c Q(\mathbf{c}|\mathbf{x}, \Theta_f) \log \left(\frac{P(\mathbf{c})}{Q(\mathbf{c}|\mathbf{x}, \Theta_f)} \right) d\mathbf{c} + \int_c Q(\mathbf{c}|\mathbf{x}, \Theta_f) \log P(\mathbf{x}|\mathbf{c}, \Theta_g) d\mathbf{c} \\ &= -D_{KL}(Q(\mathbf{c}|\mathbf{x}, \Theta_f) \| P(\mathbf{c})) + E_{(\mathbf{c}|\mathbf{x}, \Theta_f) \sim Q} [\log P(\mathbf{x}|\mathbf{c}, \Theta_g)] \end{aligned}$$

- For $P(\mathbf{c}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$:
- To minimize the first term, the encoder $f(\mathbf{x}; \Theta_f) = (\mu, \sigma)$ has a loss term [15]

$$\exp(\sigma) - (1 + \sigma) + \|\mu\|^2$$

- To maximize the second term, the decoder
 - ① Samples \mathbf{c} from $\mathcal{N}(\mu, \sigma)$ to get $g(\mathbf{c}; \Theta_g)$, the mean of output \mathcal{N}
 - ② Minimizes a loss term $\|\mathbf{x} - \mathbf{a}^{(L)}\|^2$ as in AE

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

Problems of VAE

- Only maximize a lower bound of the likelihood $P(\mathbb{X} | \Theta_f, \Theta_g)$
 - Θ_f and Θ_g are encoder and decoder weights, respectively
- Still blurry images
 - The decoder's loss is the same with that of AE

Flow-based Models

- Idea: let the decoder $g(\cdot; \Theta_g)$ be a **deterministic invertible** function
 - Given $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, we have $\mathbf{x} = g(\mathbf{c}; \Theta_g)$ of complex distribution
 - Conversely, given \mathbf{x} , we have $\mathbf{c} = g^{-1}(\mathbf{x}; \Theta_g)$
- The likelihood can be maximize directly:

$$\begin{aligned} & \arg \max_g \log P_g(\mathbb{X}) \\ &= \arg \max_g \sum_i \log P_{\mathbf{x}}(g(\mathbf{c}^{(i)})) \\ &= \arg \max_g \sum_i \log [P_{\mathbf{c}}(g^{-1}(\mathbf{x}^{(i)})) |\det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)}))|] \\ &= \arg \max_g \sum_i \log P_{\mathbf{c}}(g^{-1}(\mathbf{x}^{(i)})) + \log |\det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)}))| \end{aligned}$$

- First term: finds $g(\Theta_g)$ that maps all $\mathbf{x}^{(i)}$ to 0
- Second term: prevents $g(\Theta_g)$ from mapping all $\mathbf{x}^{(i)}$ to 0
 - $\log |\det(\mathbf{O})| = -\infty$

Flow-based Models

- Idea: let the decoder $g(\cdot; \Theta_g)$ be a **deterministic invertible** function
 - Given $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, we have $\mathbf{x} = g(\mathbf{c}; \Theta_g)$ of complex distribution
 - Conversely, given \mathbf{x} , we have $\mathbf{c} = g^{-1}(\mathbf{x}; \Theta_g)$
- The likelihood can be maximize directly:

$$\begin{aligned} & \arg \max_g \log P_g(\mathbb{X}) \\ &= \arg \max_g \sum_i \log P_{\mathbf{x}}(g(\mathbf{c}^{(i)})) \\ &= \arg \max_g \sum_i \log [P_{\mathbf{c}}(g^{-1}(\mathbf{x}^{(i)})) |\det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)}))|] \\ &= \arg \max_g \sum_i \log P_{\mathbf{c}}(g^{-1}(\mathbf{x}^{(i)})) + \log |\det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)}))| \end{aligned}$$

- First term: finds $g(\Theta_g)$ that maps all $\mathbf{x}^{(i)}$ to 0
- Second term: prevents $g(\Theta_g)$ from mapping all $\mathbf{x}^{(i)}$ to 0
 - $\log |\det(\mathbf{O})| = -\infty$
- But how to ensure the followings during training?
 - g is invertible
 - $\det(\mathbf{J}(g^{-1})(\cdot))$ can be easily computed

Ensuring Invertibility

- Glow [16]: make g^{-1} an 1×1 convolution layer

$$\bullet \quad \mathbf{x}_{i,j,:} = \begin{bmatrix} x_{i,j,1} \\ x_{i,j,2} \\ x_{i,j,3} \end{bmatrix}, \quad g^{-1} = \mathbf{W}_{3 \times 3},$$

$$g^{-1}(\mathbf{x}_{i,j,:}) = \mathbf{W}_{3 \times 3} \begin{bmatrix} x_{i,j,1} \\ x_{i,j,2} \\ x_{i,j,3} \end{bmatrix} = \begin{bmatrix} c_{i,j,1} \\ c_{i,j,2} \\ c_{i,j,3} \end{bmatrix}$$

- At training time, initialize $\mathbf{W}_{3 \times 3}$ as an invertible matrix
 - $g^{-1} = \mathbf{W}_{3 \times 3}$ is likely to be invertible after SGD updates
 - Determinant is easy to compute: $\det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)})) = \det(\mathbf{W}_{3 \times 3})^{W \times H}$
- At inference time, use $g = \mathbf{W}_{3 \times 3}^{-1}$ to generate images

Ensuring Invertibility

- Glow [16]: make g^{-1} an 1×1 convolution layer

$$\bullet \quad \mathbf{x}_{i,j,:} = \begin{bmatrix} \mathbf{x}_{i,j,1} \\ \mathbf{x}_{i,j,2} \\ \mathbf{x}_{i,j,3} \end{bmatrix}, \quad g^{-1} = \mathbf{W}_{3 \times 3},$$

$$g^{-1}(\mathbf{x}_{i,j,:}) = \mathbf{W}_{3 \times 3} \begin{bmatrix} \mathbf{x}_{i,j,1} \\ \mathbf{x}_{i,j,2} \\ \mathbf{x}_{i,j,3} \end{bmatrix} = \begin{bmatrix} \mathbf{c}_{i,j,1} \\ \mathbf{c}_{i,j,2} \\ \mathbf{c}_{i,j,3} \end{bmatrix}$$

- At training time, initialize $\mathbf{W}_{3 \times 3}$ as an invertible matrix
 - $g^{-1} = \mathbf{W}_{3 \times 3}$ is likely to be invertible after SGD updates
 - Determinant is easy to compute: $\det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)})) = \det(\mathbf{W}_{3 \times 3})^{W \times H}$
- At inference time, use $g = \mathbf{W}_{3 \times 3}^{-1}$ to generate images
- Problem: g has limited expressiveness

Step-wise Generation

- Cascade multiple invertible g to have a more complex one:

$$\mathbf{x} = g^{(T)}(\cdots g^{(2)}(g^{(1)}(\mathbf{c})))$$

Step-wise Generation

- Cascade multiple invertible g to have a more complex one:

$$\mathbf{x} = g^{(T)}(\cdots g^{(2)}(g^{(1)}(\mathbf{c})))$$

- Result: sharp images



Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

Problems of Flow-based Models

- Limited model expressiveness
- Slow training and inference
 - #steps can be large

Denoising Diffusion Probabilistic Models (DDPM) [10]

- Borrow some good ideas from previous works
 - Probabilistic formulation of VAE that models encoder in the objective
 - Step-wise encoding/decoding in generative flows
- But, unlike flows, the encoding steps
 - Are **predefined**; no parameter to learn
 - Can be simplified to **one encoding step**

Encoding

- Predefined encoding functions $f^{(t)}(\cdot; \beta^{(t)})$, $t = 1, \dots, T$:
 - $\beta^{(t)}$, $\forall t$, are hyperparameters; no learning needed
- Let each sample $\mathbf{x} = \mathbf{x}^{(0)}$
- $\mathbf{x}^{(1)} = f^{(1)}(\mathbf{x}^{(0)}; \beta^{(1)}) = \sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)} + \sqrt{\beta^{(1)}} \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - So, $(\mathbf{x}^{(1)} | \mathbf{x}^{(0)} = \mathbf{x}^{(0)}) \sim \mathcal{N}(\sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)}, \beta^{(1)} \mathbf{I})$

Encoding

- Predefined encoding functions $f^{(t)}(\cdot; \beta^{(t)})$, $t = 1, \dots, T$:
 - $\beta^{(t)}$, $\forall t$, are hyperparameters; no learning needed
- Let each sample $\mathbf{x} = \mathbf{x}^{(0)}$
- $\mathbf{x}^{(1)} = f^{(1)}(\mathbf{x}^{(0)}; \beta^{(1)}) = \sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)} + \sqrt{\beta^{(1)}} \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - So, $(\mathbf{x}^{(1)} | \mathbf{x}^{(0)} = \mathbf{x}^{(0)}) \sim \mathcal{N}(\sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)}, \beta^{(1)} \mathbf{I})$
- $$\begin{aligned}\mathbf{x}^{(2)} &= f^{(2)}(\mathbf{x}^{(1)}; \beta^{(2)}) = \sqrt{1 - \beta^{(2)}} \mathbf{x}^{(1)} + \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} \\ &= \sqrt{1 - \beta^{(2)}} \sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)} + \sqrt{1 - (1 - \beta^{(2)})(1 - \beta^{(1)})} \boldsymbol{\varepsilon} \\ &= \sqrt{\alpha^{(2)} \alpha^{(1)}} \mathbf{x}^{(0)} + \sqrt{1 - \alpha^{(2)} \alpha^{(1)}} \boldsymbol{\varepsilon}\end{aligned}$$
 - $\sqrt{1 - \beta^{(1)}} \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} + \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, ((1 - \beta^{(1)}) \beta^{(2)} + \beta^{(2)}) \mathbf{I})$
 - Let $\alpha^{(t)} = 1 - \beta^{(t)}$ (derived hyperparameter)
 - Only one $\boldsymbol{\varepsilon}$ is added

Encoding

- Predefined encoding functions $f^{(t)}(\cdot; \beta^{(t)})$, $t = 1, \dots, T$:
 - $\beta^{(t)}$, $\forall t$, are hyperparameters; no learning needed
- Let each sample $\mathbf{x} = \mathbf{x}^{(0)}$
- $\mathbf{x}^{(1)} = f^{(1)}(\mathbf{x}^{(0)}; \beta^{(1)}) = \sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)} + \sqrt{\beta^{(1)}} \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - So, $(\mathbf{x}^{(1)} | \mathbf{x}^{(0)} = \mathbf{x}^{(0)}) \sim \mathcal{N}(\sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)}, \beta^{(1)} \mathbf{I})$
- $$\begin{aligned}\mathbf{x}^{(2)} &= f^{(2)}(\mathbf{x}^{(1)}; \beta^{(2)}) = \sqrt{1 - \beta^{(2)}} \mathbf{x}^{(1)} + \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} \\ &= \sqrt{1 - \beta^{(2)}} \sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)} + \sqrt{1 - (1 - \beta^{(2)})(1 - \beta^{(1)})} \boldsymbol{\varepsilon} \\ &= \sqrt{\alpha^{(2)} \alpha^{(1)}} \mathbf{x}^{(0)} + \sqrt{1 - \alpha^{(2)} \alpha^{(1)}} \boldsymbol{\varepsilon}\end{aligned}$$
 - $\sqrt{1 - \beta^{(1)}} \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} + \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, ((1 - \beta^{(1)}) \beta^{(2)} + \beta^{(2)}) \mathbf{I})$
 - Let $\bar{\alpha}^{(t)} = 1 - \beta^{(t)}$ (derived hyperparameter)
 - Only one $\boldsymbol{\varepsilon}$ is added
- $\mathbf{x}^{(t)} = \sqrt{\bar{\alpha}^{(t)}} \mathbf{x}^{(t-1)} + \sqrt{1 - \bar{\alpha}^{(t)}} \boldsymbol{\varepsilon}$
 - Let $\bar{\alpha}^{(t)} = \alpha^{(t)} \alpha^{(t-1)} \dots \alpha^{(1)}$ (derived hyperparameter)

Objective

- As VAE, DDPM maximizes a lower bound of $P(\mathbb{X})$
- VAE: for each \mathbf{x} , maximize

$$\int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{x}, \mathbf{c})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} = E_{\mathbf{c}|\mathbf{x} \sim Q} \left[\log \left(\frac{P(\mathbf{x}, \mathbf{c})}{Q(\mathbf{c}|\mathbf{x})} \right) \right] \leq P(\mathbf{x})$$

- DDPM: for each \mathbf{x} , maximize

$$E_{(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)} | \mathbf{x}^{(0)}) \sim Q} \left[\log \left(\frac{P(\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)})}{Q(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)} | \mathbf{x}^{(0)})} \right) \right],$$

which can be simplified to [22]:

$$-D_{KL}(Q(\mathbf{x}^{(T)} | \mathbf{x}^{(0)}) \| P(\mathbf{x}^{(T)})) + E_{\mathbf{x}^{(1)} | \mathbf{x}^{(0)} \sim Q} [\log P(\mathbf{x}^{(0)} | \mathbf{x}^{(1)})] + \\ - \sum_{t=2}^T E_{\mathbf{x}^{(t)} | \mathbf{x}^{(0)} \sim Q} [D_{KL}(Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \| P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}))]$$

- First term controlled by encoding process (predefined)
- Second & third term controlled by decoding process (learnable)

Decoding I

- For simplicity, we focus on maximizing the third term:

$$-\sum_{t=2}^T E_{\mathbf{x}^{(t)} | \mathbf{x}^{(0)} \sim Q} \left[D_{KL} \left(Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \| P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \right) \right]$$

- Goal: for each observed $\mathbf{x}^{(t)}$, minimize

$$D_{KL} \left(Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \| P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \right)$$

- Note that

$$\begin{aligned} Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) &= \frac{Q(\mathbf{x}^{(t-1)}, \mathbf{x}^{(t)}, \mathbf{x}^{(0)})}{Q(\mathbf{x}^{(t)}, \mathbf{x}^{(0)})} \\ &= \frac{Q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}) Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(0)}) Q(\mathbf{x}^{(0)})}{Q(\mathbf{x}^{(t)} | \mathbf{x}^{(0)}) Q(\mathbf{x}^{(0)})} \\ &= \frac{Q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}) Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(0)})}{Q(\mathbf{x}^{(t)} | \mathbf{x}^{(0)})} \end{aligned}$$

- Since $Q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$ & $Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(0)})$ are Gaussian, we have [22]:

$$Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) = \mathcal{N} \left(\frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}}, \frac{1-\bar{\alpha}^{(t-1)}}{1-\bar{\alpha}^{(t)}}\beta^{(t)}\mathbf{I} \right)$$

Decoding II

- Goal: for each observed $\mathbf{x}^{(t)}$, minimize

$$D_{KL} \left(Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \| P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \right),$$

where $Q(\dots) = \mathcal{N} \left(\frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}}, \dots \right)$ is **fixed**

- DDPM finds Θ that move the mean of $P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \Theta)$ (also Gaussian) toward $Q(\dots)$'s mean:

$$\frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}}$$

Noise Predictor

- Θ moves $P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \Theta)$'s mean toward

$$\begin{aligned}& \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}} \\&= \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\frac{\mathbf{x}^{(t)} - \sqrt{1-\bar{\alpha}^{(t)}}\boldsymbol{\varepsilon}}{\sqrt{\bar{\alpha}^{(t)}}}}{1-\bar{\alpha}^{(t)}} \\&= \frac{1}{\sqrt{\alpha^{(t)}}} \left(\mathbf{x}^{(t)} - \frac{1-\alpha^{(t)}}{\sqrt{1-\bar{\alpha}^{(t)}}} \boldsymbol{\varepsilon} \right)\end{aligned}$$

- What's its corresponding network?

Noise Predictor

- Θ moves $P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \Theta)$'s mean toward

$$\begin{aligned}& \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}} \\&= \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\frac{\mathbf{x}^{(t)} - \sqrt{1-\bar{\alpha}^{(t)}}\boldsymbol{\varepsilon}}{\sqrt{\bar{\alpha}^{(t)}}}}{1-\bar{\alpha}^{(t)}} \\&= \frac{1}{\sqrt{\alpha^{(t)}}} \left(\mathbf{x}^{(t)} - \frac{1-\alpha^{(t)}}{\sqrt{1-\bar{\alpha}^{(t)}}} \boldsymbol{\varepsilon} \right)\end{aligned}$$

- What's its corresponding network?
- By definition: let Θ parametrize a network outputting $\mathbf{x}^{(t-1)}$ given $\mathbf{x}^{(t)}$
 - But the input $\mathbf{x}^{(t)}$ also resides in the output target

Noise Predictor

- Θ moves $P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \Theta)$'s mean toward

$$\begin{aligned}& \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}} \\&= \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\frac{\mathbf{x}^{(t)} - \sqrt{1-\bar{\alpha}^{(t)}}\boldsymbol{\varepsilon}}{\sqrt{\bar{\alpha}^{(t)}}}}{1-\bar{\alpha}^{(t)}} \\&= \frac{1}{\sqrt{\alpha^{(t)}}} \left(\mathbf{x}^{(t)} - \frac{1-\alpha^{(t)}}{\sqrt{1-\bar{\alpha}^{(t)}}} \boldsymbol{\varepsilon} \right)\end{aligned}$$

- What's its corresponding network?
- By definition: let Θ parametrize a network outputting $\mathbf{x}^{(t-1)}$ given $\mathbf{x}^{(t)}$
 - But the input $\mathbf{x}^{(t)}$ also resides in the output target
- DDPM: let Θ parametrize a **noise predictor** outputting $\boldsymbol{\varepsilon}$ given $\mathbf{x}^{(t)}$
 - Objective: $\arg \min_{\Theta} \|\boldsymbol{\varepsilon} - e(\mathbf{x}^{(t)}, t; \Theta)\|^2$
 - Shared between all t

Training & Inference Algorithms

- One-step encoding during training time
- Multi-step inference (sampling), with each intermediate decoding step t , $t > 1$, comes with extra noise $\sigma^{(t)} \mathbf{z}$
 - Similar to output token sampling in GPT
 - Improves performance empirically

Algorithm 1 Training

```
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
     
$$\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$$

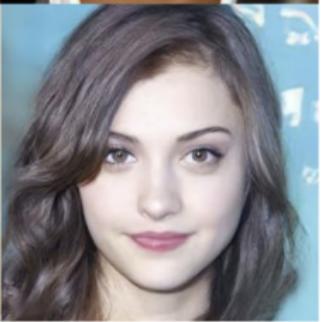
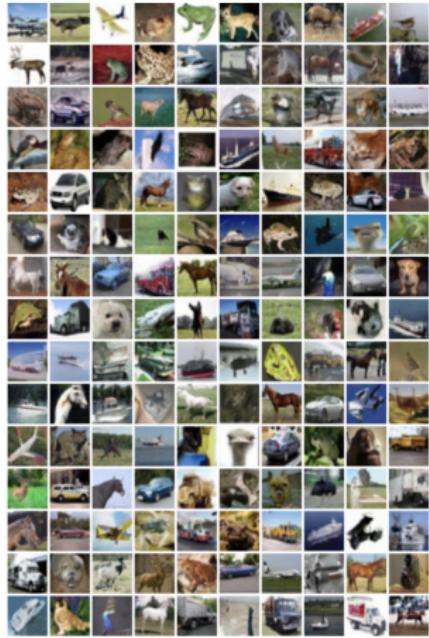
6: until converged
```

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

Results

- Sharp and coherent



Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

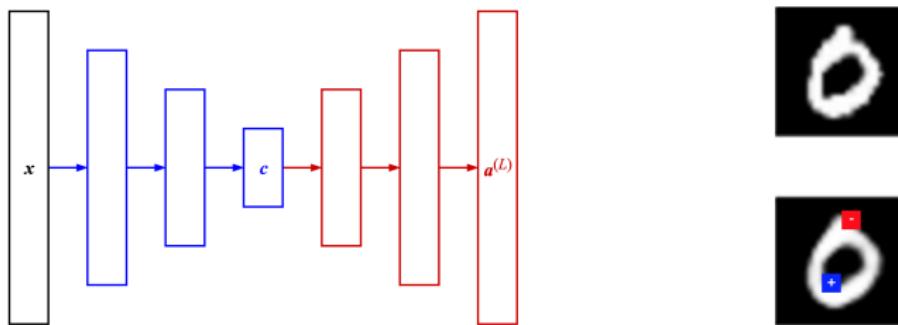
⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

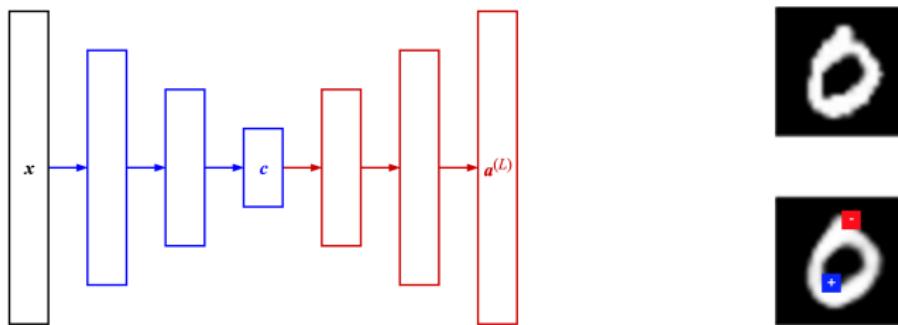
Why Do AE and VAE Give Blurry Images?

- Objective $\arg \max_{\Theta} \log P(\mathbb{X} | \Theta) = \arg \min_{\Theta} \sum_n \|x^{(n)} - a^{(n,L)}\|^2$ does not penalize Gaussian pixel noises in $a^{(i,L)}$
 - Root cause: assumption that $x \sim \mathcal{N}$



Why Do AE and VAE Give Blurry Images?

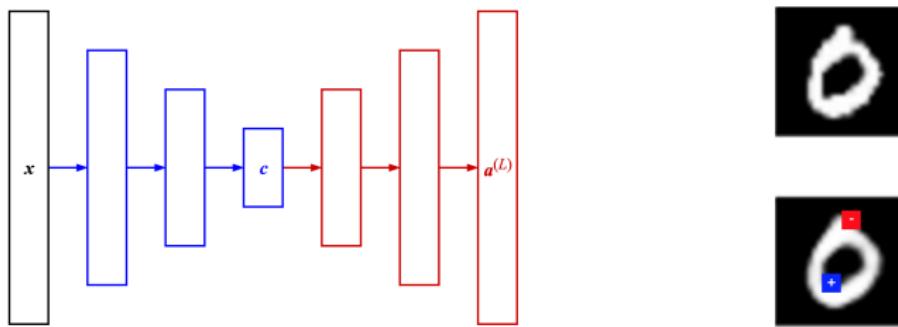
- Objective $\arg \max_{\Theta} \log P(\mathbb{X} | \Theta) = \arg \min_{\Theta} \sum_n \|x^{(n)} - a^{(n,L)}\|^2$ does not penalize Gaussian pixel noises in $a^{(i,L)}$
 - Root cause: assumption that $x \sim \mathcal{N}$



- Fix 1: maximizing the likelihood $P(\mathbb{X} | \Theta)$ without assuming $x \sim \mathcal{N}$
 - Flow-based methods, diffusion models

Why Do AE and VAE Give Blurry Images?

- Objective $\arg \max_{\Theta} \log P(\mathbb{X} | \Theta) = \arg \min_{\Theta} \sum_n \|x^{(n)} - a^{(n,L)}\|^2$ does not penalize Gaussian pixel noises in $a^{(i,L)}$
 - Root cause: assumption that $x \sim \mathcal{N}$



- Fix 1: maximizing the likelihood $P(\mathbb{X} | \Theta)$ without assuming $x \sim \mathcal{N}$
 - Flow-based methods, diffusion models
- Fix 2: maximizing the chance that $a^{(L)}$ *is an true image from another NN point of view*

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

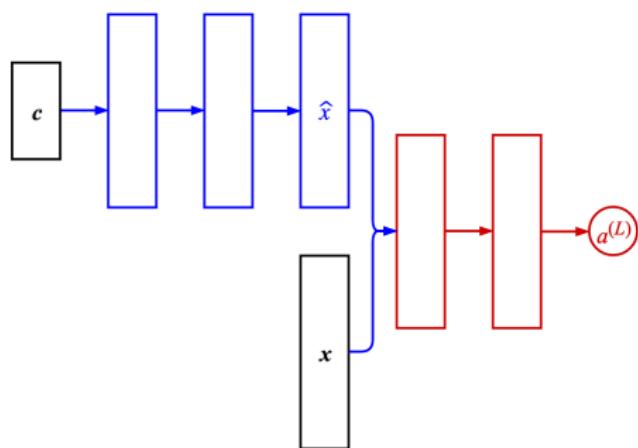
⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

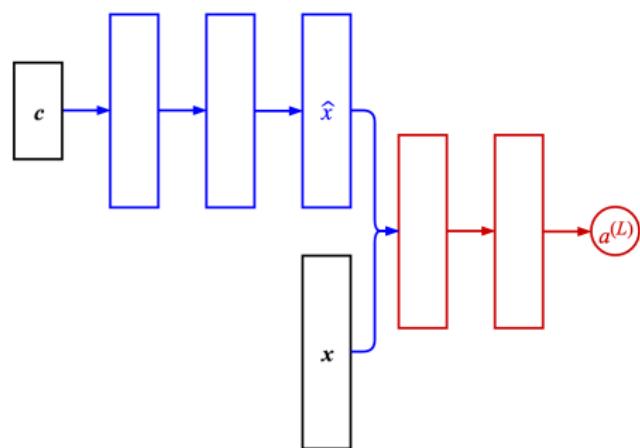
Network Design

- *Generative adversarial networks (GAN)* [8] consist of:



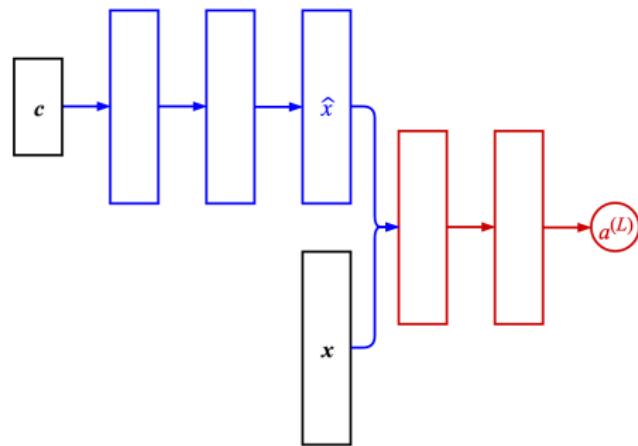
Network Design

- **Generative adversarial networks (GAN)** [8] consist of:
- **Generator g** : to generate data points from random codes
 - No need for “encoder” since the task is data synthesis



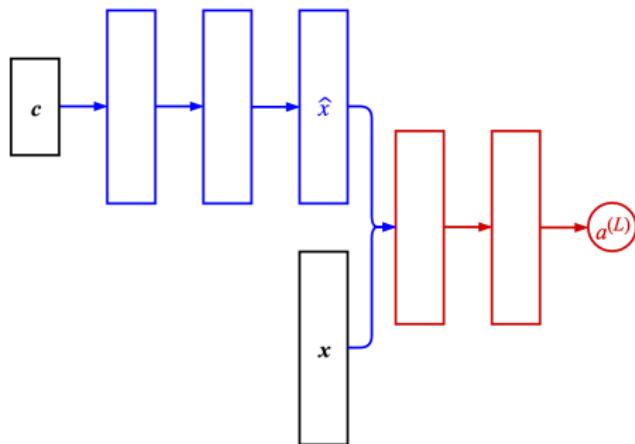
Network Design

- **Generative adversarial networks (GAN)** [8] consist of:
 - **Generator g** : to generate data points from random codes
 - No need for “encoder” since the task is data synthesis
 - **Discriminator f** : to separate generated points from real ones
 - Weights for x and \hat{x} are tied
 - A binary classifier with Sigmoid output unit $a^{(L)} = \hat{p}$ for $P(y = \text{true point} | \mathbf{x}) \sim \text{Bernoulli}(\rho)$



Network Design

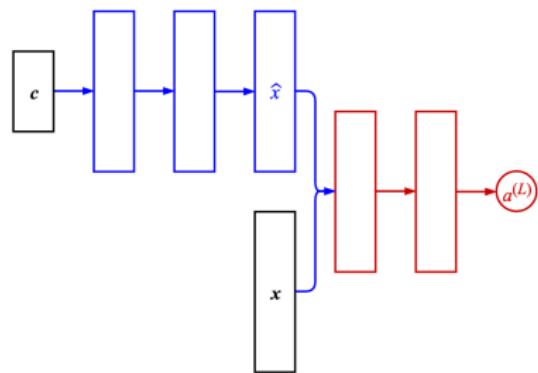
- **Generative adversarial networks (GAN)** [8] consist of:
 - **Generator g** : to generate data points from random codes
 - No need for “encoder” since the task is data synthesis
 - **Discriminator f** : to separate generated points from real ones
 - Weights for x and \hat{x} are tied
 - A binary classifier with Sigmoid output unit $a^{(L)} = \hat{p}$ for $P(y = \text{true point} | \mathbf{x}) \sim \text{Bernoulli}(\rho)$
- Goal: to train a g that tricks f into believing $g(c)$ is real



Cost Function

- Given N real training points and N generated points:

$$\begin{aligned} & \arg \min_{\Theta_g} \max_{\Theta_f} \log P(\mathbb{X} | \Theta_g, \Theta_f) \\ &= \arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)}))) \end{aligned}$$



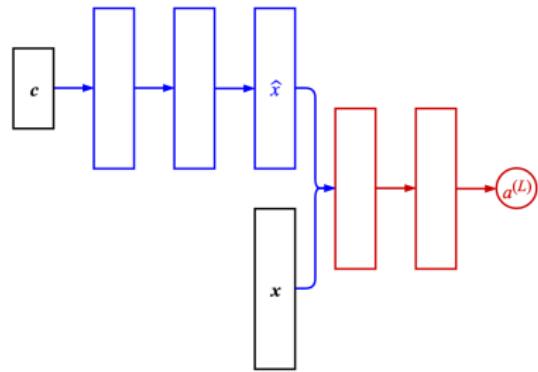
Cost Function

- Given N real training points and N generated points:

$$\begin{aligned} & \arg \min_{\Theta_g} \max_{\Theta_f} \log P(\mathbb{X} | \Theta_g, \Theta_f) \\ &= \arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)}))) \\ &= \arg \min_{\Theta_g} \max_{\Theta_f} \sum_{n=1}^N \log \hat{\rho}^{(n)} + \sum_{m=1}^N \log(1 - \hat{\rho}^{(m)}) \end{aligned}$$

- Recall that f maximizes the log likelihood

$$\log P(\mathbb{X} | \Theta) \propto \sum_n \log P(y^{(n)} | \mathbf{x}^{(n)}, \Theta) = \sum_n \log \left[(\hat{\rho}^{(n)})^{y^{(n)}} (1 - \hat{\rho}^{(n)})^{(1-y^{(n)})} \right]$$



Cost Function

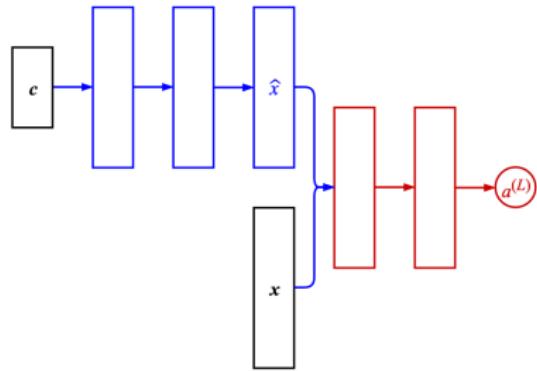
- Given N real training points and N generated points:

$$\begin{aligned} & \arg \min_{\Theta_g} \max_{\Theta_f} \log P(\mathbb{X} | \Theta_g, \Theta_f) \\ &= \arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)}))) \\ &= \arg \min_{\Theta_g} \max_{\Theta_f} \sum_{n=1}^N \log \hat{\rho}^{(n)} + \sum_{m=1}^N \log(1 - \hat{\rho}^{(m)}) \end{aligned}$$

- Recall that f maximizes the log likelihood

$$\log P(\mathbb{X} | \Theta) \propto \sum_n \log P(y^{(n)} | \mathbf{x}^{(n)}, \Theta) = \sum_n \log \left[(\hat{\rho}^{(n)})^{y^{(n)}} (1 - \hat{\rho}^{(n)})^{(1-y^{(n)})} \right]$$

- Inner **max** first, then outer **min**



Cost Function

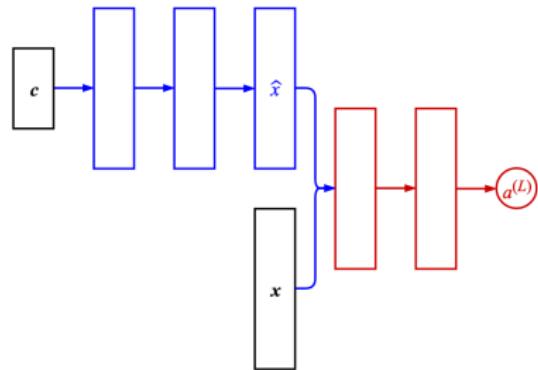
- Given N real training points and N generated points:

$$\begin{aligned} & \arg \min_{\Theta_g} \max_{\Theta_f} \log P(\mathbb{X} | \Theta_g, \Theta_f) \\ &= \arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)}))) \\ &= \arg \min_{\Theta_g} \max_{\Theta_f} \sum_{n=1}^N \log \hat{\rho}^{(n)} + \sum_{m=1}^N \log(1 - \hat{\rho}^{(m)}) \end{aligned}$$

- Recall that f maximizes the log likelihood

$$\log P(\mathbb{X} | \Theta) \propto \sum_n \log P(y^{(n)} | \mathbf{x}^{(n)}, \Theta) = \sum_n \log \left[(\hat{\rho}^{(n)})^{y^{(n)}} (1 - \hat{\rho}^{(n)})^{(1-y^{(n)})} \right]$$

- Inner **max** first, then outer **min**
- $\hat{\rho}^{(n)}$ depends on Θ_f only
- $\hat{\rho}^{(m)}$ depends on both Θ_f and Θ_g



Training: Alternative SGD

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

- Initialize Θ_g for g and Θ_f for f
- At each SGD step/iteration:

Training: Alternative SGD

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

- Initialize Θ_g for g and Θ_f for f
- At each SGD step/iteration:

① *Repeat K times (with fixed Θ_g):*

- ① Sample N real points $\{\mathbf{x}^{(n)}\}_n$ from \mathbb{X} and N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ② $\Theta_f \leftarrow \Theta_f + \eta \nabla_{\Theta_f} [\sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$

Training: Alternative SGD

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

- Initialize Θ_g for g and Θ_f for f
- At each SGD step/iteration:

① *Repeat K times (with fixed Θ_g):*

- ① Sample N real points $\{\mathbf{x}^{(n)}\}_n$ from \mathbb{X} and N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ② $\Theta_f \leftarrow \Theta_f + \eta \nabla_{\Theta_f} [\sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$

② *Execute once (with fixed Θ_f):*

- ① Sample N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ② $\Theta_g \leftarrow \Theta_g - \eta \nabla_{\Theta_g} [\sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$

Training: Alternative SGD

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

- Initialize Θ_g for g and Θ_f for f
- At each SGD step/iteration:

① *Repeat K times (with fixed Θ_g):*

- ① Sample N real points $\{\mathbf{x}^{(n)}\}_n$ from \mathbb{X} and N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ② $\Theta_f \leftarrow \Theta_f + \eta \nabla_{\Theta_f} [\sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$

② *Execute once (with fixed Θ_f):*

- ① Sample N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ② $\Theta_g \leftarrow \Theta_g - \eta \nabla_{\Theta_g} [\sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$

- Why limiting the steps (K) when updating Θ_f ?

Training: Alternative SGD

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

- Initialize Θ_g for g and Θ_f for f
- At each SGD step/iteration:

① ***Repeat K times (with fixed Θ_g):***

- ① Sample N real points $\{\mathbf{x}^{(n)}\}_n$ from \mathbb{X} and N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ② $\Theta_f \leftarrow \Theta_f + \eta \nabla_{\Theta_f} [\sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$

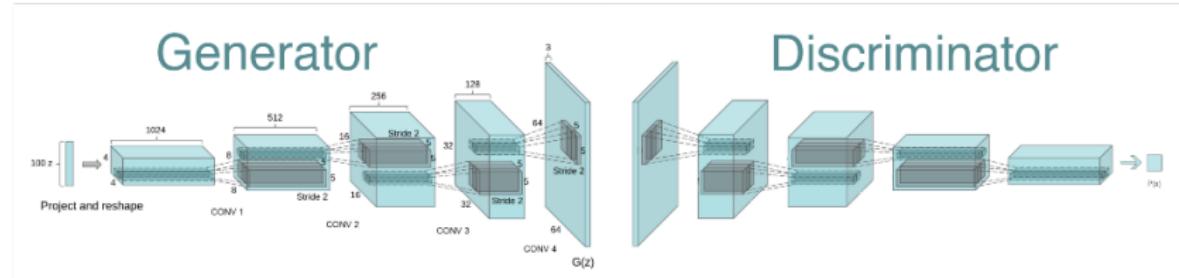
② ***Execute once (with fixed Θ_f):***

- ① Sample N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ② $\Theta_g \leftarrow \Theta_g - \eta \nabla_{\Theta_g} [\sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$

- Why limiting the steps (K) when updating Θ_f ?
 - f may overfit data and give very different values once g is updated
 - Limiting K so to prevent g from being updated for “wrong” target

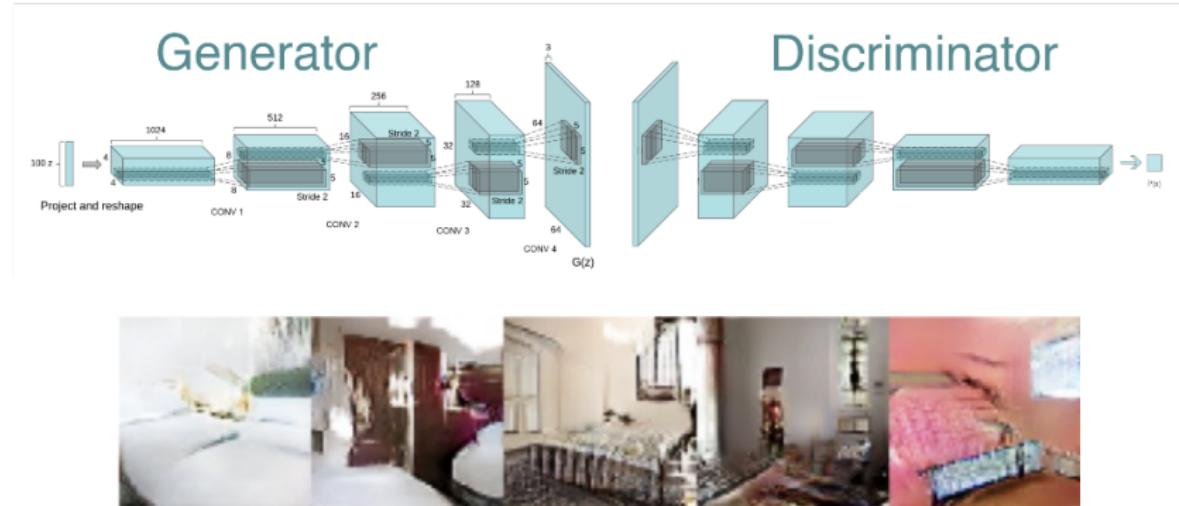
Results

- Domain-specific architecture, e.g., DC-GAN [29]



Results

- Domain-specific architecture, e.g., DC-GAN [29]



GANs Are Hard to Train!

- [Tips for Training Stable GANs](#)
- [Keep Calm and train a GAN. Pitfalls and Tips...](#)
- [10 Lessons I Learned Training GANs for one Year](#)
- [GAN hacks on GitHub](#)

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

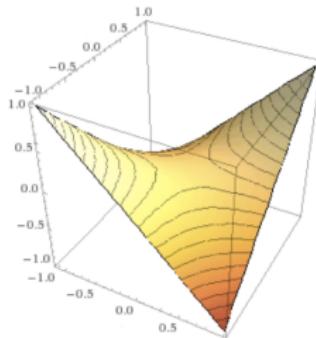
Challenge: Non-Convergence

- The GAN training may *not* converge

Challenge: Non-Convergence

- The GAN training may *not* converge
- The goal of GAN is to find a saddle point

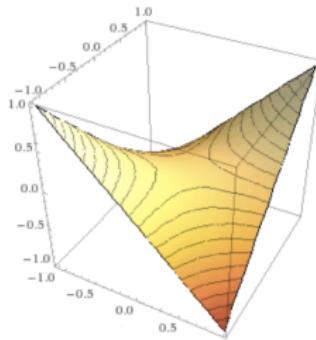
$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$



Challenge: Non-Convergence

- The GAN training may *not* converge
- The goal of GAN is to find a saddle point

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

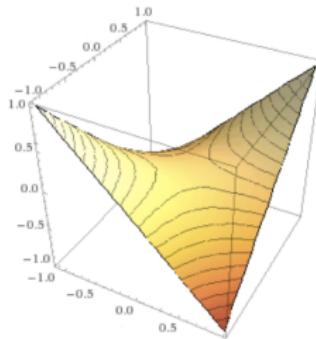


- The updated Θ_f and Θ_g may cancel each other's progress

Challenge: Non-Convergence

- The GAN training may *not* converge
- The goal of GAN is to find a saddle point

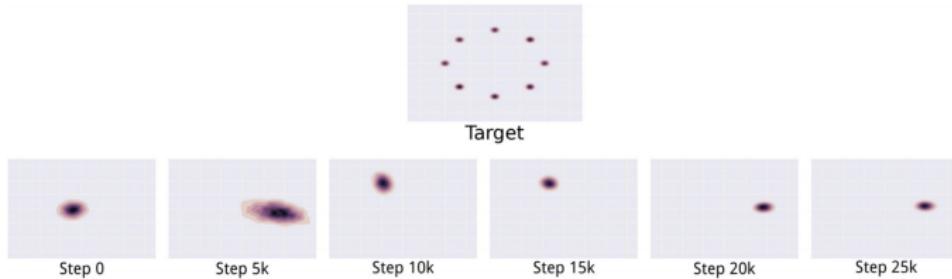
$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$



- The updated Θ_f and Θ_g may cancel each other's progress
- Requires human monitoring and termination

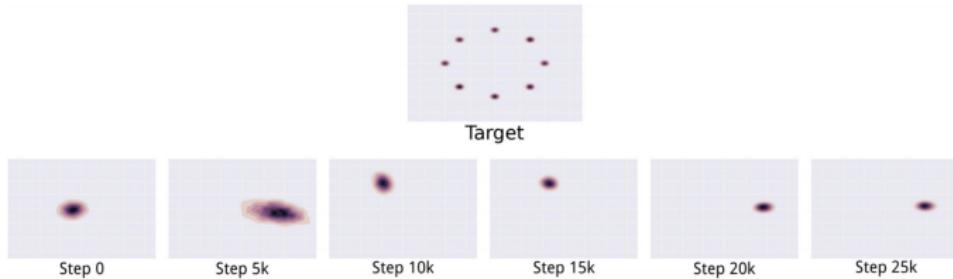
Mode Collapsing

- Even worse: *mode collapsing*
 - g may oscillate from generating one kind of points to generating another kind of points



Mode Collapsing

- Even worse: *mode collapsing*
 - g may oscillate from generating one kind of points to generating another kind of points



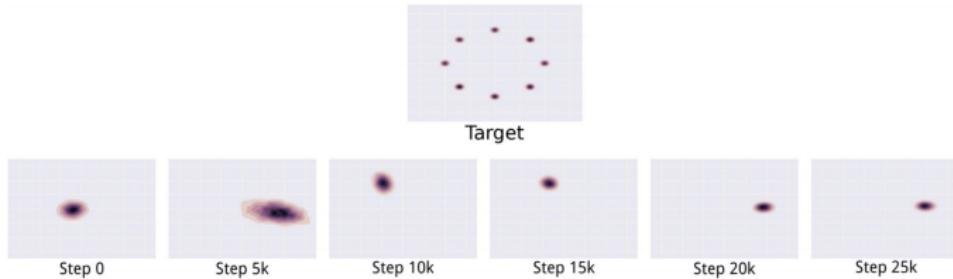
- When K is small, alternate SGD does not distinguish between $\min_{\Theta_g} \max_{\Theta_f}$ and $\max_{\Theta_f} \min_{\Theta_g}$

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

- $\max_{\Theta_f} \min_{\Theta_g}$?

Mode Collapsing

- Even worse: *mode collapsing*
 - g may oscillate from generating one kind of points to generating another kind of points



- When K is small, alternate SGD does not distinguish between $\min_{\Theta_g} \max_{\Theta_f}$ and $\max_{\Theta_f} \min_{\Theta_g}$
- $$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$
- $\max_{\Theta_f} \min_{\Theta_g}$? g is encouraged to map every code to the “mode” that f believes is most likely to be real

Solutions

- Minibatch discrimination [35]
 - In $\max_{\Theta_f} \min_{\Theta_g}$ case, g collapses because $\nabla_{\Theta_f} C$ are computed independently for each point
 - Why not augment each $x^{(n)}/\hat{x}^{(n)}$ with **batch features**?

Solutions

- Minibatch discrimination [35]
 - In $\max_{\Theta_f} \min_{\Theta_g}$ case, g collapses because $\nabla_{\Theta_f} C$ are computed independently for each point
 - Why not augment each $x^{(n)}/\hat{x}^{(n)}$ with **batch features**?
 - If g collapses, f can tell this from batch features and reject fake points
 - Now, g needs to generate dissimilar points to fool f

Solutions

- Minibatch discrimination [35]

- In $\max_{\Theta_f} \min_{\Theta_g}$ case, g collapses because $\nabla_{\Theta_f} C$ are computed independently for each point
- Why not augment each $x^{(n)}/\hat{x}^{(n)}$ with **batch features**?
- If g collapses, f can tell this from batch features and reject fake points
- Now, g needs to generate dissimilar points to fool f

6	0	6	3	7	3	4	5	2	8
1	0	1	0	7	9	7	1	4	0
7	6	5	3	8	3	6	9	9	6
9	8	4	6	2	3	=	0	8	7
0	1	5	8	4	1	3	7	6	9
1	4	2	3	7	1	1	4	8	5
2	8	1	4	3	5	1	1	2	4
1	7	0	4	3	6	2	4	7	9
6	8	3	1	6	7	5	0	7	6
9	7	7	0	8	1	3	1	6	9

without

5	8	9	9	9	5	0	2	8	1
4	6	6	6	0	7	1	2	7	8
2	0	1	0	7	4	4	9	5	9
1	7	3	8	3	2	1	6	3	8
0	2	6	0	6	9	3	6	1	6
5	7	9	9	5	8	6	4	4	5
1	5	4	6	7	8	7	9	7	3
7	4	4	7	9	5	6	8	1	1
1	3	9	0	8	1	1	1	8	7
2	6	8	3	0	9	7	1	4	1

with

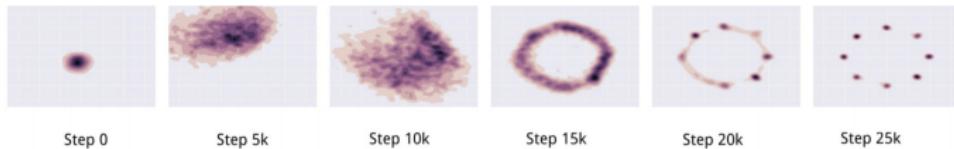
Solutions

- Minibatch discrimination [35]

- In $\max_{\Theta_f} \min_{\Theta_g}$ case, g collapses because $\nabla_{\Theta_f} C$ are computed independently for each point
- Why not augment each $x^{(n)}/\hat{x}^{(n)}$ with **batch features**?
- If g collapses, f can tell this from batch features and reject fake points
- Now, g needs to generate dissimilar points to fool f

	<table border="1"><tr><td>6</td><td>3</td><td>6</td><td>7</td><td>3</td><td>4</td><td>5</td><td>2</td><td>8</td><td>1</td></tr><tr><td>1</td><td>0</td><td>1</td><td>0</td><td>7</td><td>9</td><td>7</td><td>1</td><td>4</td><td>0</td></tr><tr><td>7</td><td>6</td><td>5</td><td>3</td><td>8</td><td>3</td><td>6</td><td>9</td><td>9</td><td>0</td></tr><tr><td>4</td><td>8</td><td>4</td><td>6</td><td>2</td><td>3</td><td>=</td><td>9</td><td>8</td><td>7</td></tr><tr><td>0</td><td>1</td><td>5</td><td>8</td><td>4</td><td>1</td><td>3</td><td>7</td><td>6</td><td>9</td></tr><tr><td>1</td><td>4</td><td>2</td><td>3</td><td>7</td><td>1</td><td>1</td><td>4</td><td>8</td><td>5</td></tr><tr><td>2</td><td>8</td><td>1</td><td>4</td><td>3</td><td>5</td><td>1</td><td>1</td><td>2</td><td>9</td></tr><tr><td>1</td><td>7</td><td>0</td><td>4</td><td>3</td><td>6</td><td>2</td><td>4</td><td>7</td><td>9</td></tr><tr><td>6</td><td>8</td><td>3</td><td>1</td><td>6</td><td>7</td><td>5</td><td>0</td><td>7</td><td>6</td></tr><tr><td>9</td><td>7</td><td>2</td><td>0</td><td>3</td><td>1</td><td>3</td><td>1</td><td>6</td><td>9</td></tr></table>	6	3	6	7	3	4	5	2	8	1	1	0	1	0	7	9	7	1	4	0	7	6	5	3	8	3	6	9	9	0	4	8	4	6	2	3	=	9	8	7	0	1	5	8	4	1	3	7	6	9	1	4	2	3	7	1	1	4	8	5	2	8	1	4	3	5	1	1	2	9	1	7	0	4	3	6	2	4	7	9	6	8	3	1	6	7	5	0	7	6	9	7	2	0	3	1	3	1	6	9
6	3	6	7	3	4	5	2	8	1																																																																																												
1	0	1	0	7	9	7	1	4	0																																																																																												
7	6	5	3	8	3	6	9	9	0																																																																																												
4	8	4	6	2	3	=	9	8	7																																																																																												
0	1	5	8	4	1	3	7	6	9																																																																																												
1	4	2	3	7	1	1	4	8	5																																																																																												
2	8	1	4	3	5	1	1	2	9																																																																																												
1	7	0	4	3	6	2	4	7	9																																																																																												
6	8	3	1	6	7	5	0	7	6																																																																																												
9	7	2	0	3	1	3	1	6	9																																																																																												

- Unrolled GANs [23]: to back-propagate through **several max steps** when computing $\nabla_{\Theta_g} C$



Challenge: Balance between g and f

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

- Alternate SGD:

- $\Theta_f \leftarrow \Theta_f + \eta \nabla_{\Theta_f} [\sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$ for K times
- $\Theta_g \leftarrow \Theta_g - \eta \nabla_{\Theta_g} [\sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$

Challenge: Balance between g and f

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

- Alternate SGD:
 - $\Theta_f \leftarrow \Theta_f + \eta \nabla_{\Theta_f} [\sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$ for K times
 - $\Theta_g \leftarrow \Theta_g - \eta \nabla_{\Theta_g} [\sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$
- Why limiting K when updating Θ_f ?
- Too large K :
- Too small K :

Challenge: Balance between g and f

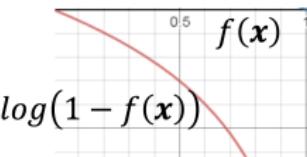
$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

- Alternate SGD:
 - $\Theta_f \leftarrow \Theta_f + \eta \nabla_{\Theta_f} [\sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$ for K times
 - $\Theta_g \leftarrow \Theta_g - \eta \nabla_{\Theta_g} [\sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$
- Why limiting K when updating Θ_f ?
- Too large K :
 - f may overfit data, making g updated for “wrong” target f
- Too small K :

Challenge: Balance between g and f

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

- Alternate SGD:
 - $\Theta_f \leftarrow \Theta_f + \eta \nabla_{\Theta_f} [\sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$ for K times
 - $\Theta_g \leftarrow \Theta_g - \eta \nabla_{\Theta_g} [\sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$
- Why limiting K when updating Θ_f ?
- Too large K :
 - f may overfit data, making g updated for “wrong” target f
 - Vanishing gradients: $\nabla_{\Theta_g} [\sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$ too small to learn

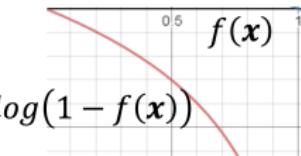


- Too small K :

Challenge: Balance between g and f

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

- Alternate SGD:
 - $\Theta_f \leftarrow \Theta_f + \eta \nabla_{\Theta_f} [\sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$ for K times
 - $\Theta_g \leftarrow \Theta_g - \eta \nabla_{\Theta_g} [\sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$
- Why limiting K when updating Θ_f ?
- Too large K :
 - f may overfit data, making g updated for “wrong” target f
 - Vanishing gradients: $\nabla_{\Theta_g} [\sum_m \log(1 - f(g(\mathbf{c}^{(m)})))]$ too small to learn



- Too small K :
 - g updated for “meaningless” f

Solution: Wasserstein GAN [1]

- Let f be a regressor **without** the sigmoid output layer
- Cost function:

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)}))$$

- Initialize Θ_g for g and Θ_f for f
- At each SGD step/iteration:

Solution: Wasserstein GAN [1]

- Let f be a regressor **without** the sigmoid output layer
- Cost function:

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)}))$$

- Initialize Θ_g for g and Θ_f for f
- At each SGD step/iteration:
 - Repeat K times (with fixed Θ_g):
 - Sample N real points $\{\mathbf{x}^{(n)}\}_n$ from \mathbb{X} and N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - $\Theta_f \leftarrow \Theta_f + \eta \text{clip}(\nabla_{\Theta_f} [\sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)}))])$

Solution: Wasserstein GAN [1]

- Let f be a regressor **without** the sigmoid output layer
- Cost function:

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)}))$$

- Initialize Θ_g for g and Θ_f for f
- At each SGD step/iteration:
 - Repeat K times (with fixed Θ_g):
 - Sample N real points $\{\mathbf{x}^{(n)}\}_n$ from \mathbb{X} and N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - $\Theta_f \leftarrow \Theta_f + \eta \text{clip}(\nabla_{\Theta_f} [\sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)}))])$
 - Execute once (with fixed Θ_f):
 - Sample N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - $\Theta_g \leftarrow \Theta_g - \eta \nabla_{\Theta_g} [-\sum_m f(g(\mathbf{c}^{(m)}))]$

GANs from Information Theory Perspective

- Review the Information Theory first!

GANs from Information Theory Perspective

- Review the Information Theory first!
- Let P_{data} / P_g be distribution of $\mathbf{x} / \hat{\mathbf{x}} = g(\mathbf{c})$

GANs from Information Theory Perspective

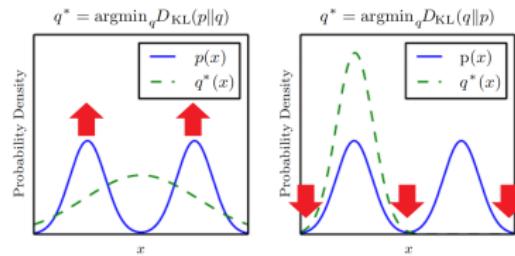
- Review the Information Theory first!
- Let P_{data} / P_g be distribution of $\mathbf{x} / \hat{\mathbf{x}} = g(\mathbf{c})$
- A way to find g : $\arg \min_{\Theta_g} D_{\text{KL}}(P_{\text{data}} \| P_g)$

GANs from Information Theory Perspective

- Review the Information Theory first!
- Let P_{data} / P_g be distribution of $\mathbf{x} / \hat{\mathbf{x}} = g(\mathbf{c})$
- A way to find g : $\arg \min_{\Theta_g} D_{\text{KL}}(P_{\text{data}} \| P_g)$
 - Why not $\arg \min_{\Theta_g} D_{\text{KL}}(P_g \| P_{\text{data}})$?

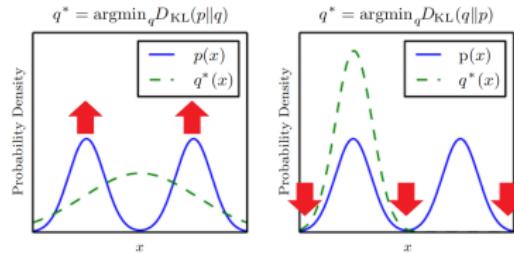
GANs from Information Theory Perspective

- Review the Information Theory first!
- Let P_{data} / P_g be distribution of $\mathbf{x} / \hat{\mathbf{x}} = g(\mathbf{c})$
- A way to find g : $\arg \min_{\Theta_g} D_{\text{KL}}(P_{\text{data}} \| P_g)$
 - Why not $\arg \min_{\Theta_g} D_{\text{KL}}(P_g \| P_{\text{data}})$?



GANs from Information Theory Perspective

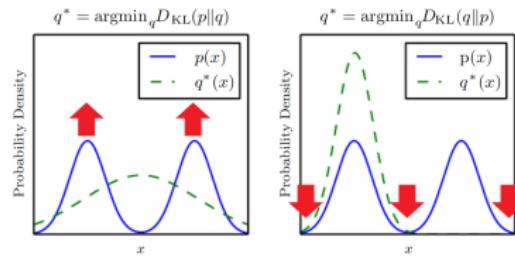
- Review the Information Theory first!
- Let P_{data} / P_g be distribution of $\mathbf{x} / \hat{\mathbf{x}} = g(\mathbf{c})$
- A way to find g : $\arg \min_{\Theta_g} D_{\text{KL}}(P_{\text{data}} \| P_g)$
 - Why not $\arg \min_{\Theta_g} D_{\text{KL}}(P_g \| P_{\text{data}})$?



- GAN: $\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(x^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$

GANs from Information Theory Perspective

- Review the Information Theory first!
- Let P_{data} / P_g be distribution of $\mathbf{x} / \hat{\mathbf{x}} = g(\mathbf{c})$
- A way to find g : $\arg \min_{\Theta_g} D_{\text{KL}}(P_{\text{data}} \| P_g)$
 - Why not $\arg \min_{\Theta_g} D_{\text{KL}}(P_g \| P_{\text{data}})$?



- GAN: $\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(x^{(n)}) + \sum_m \log(1 - f(g(c^{(m)})))$
- Actually, the max term measures Jensen-Shannon divergence (a.k.a. symmetric KL divergence):

$$D_{\text{JS}}(P_{\text{data}} \| P_g) = \frac{1}{2} D_{\text{KL}}(P_{\text{data}} \| Q) + \frac{1}{2} D_{\text{KL}}(P_g \| Q), \text{ where } Q = \frac{1}{2}(P_g + P_{\text{data}})$$

Why Jensen-Shannon Divergence? I

- Given a fixed g , we have

$$C^* = \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)})))$$

Why Jensen-Shannon Divergence? I

- Given a fixed g , we have

$$\begin{aligned} C^* &= \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)}))) \\ &= \max_{\Theta_f} \frac{1}{N} \sum_n \log f(\mathbf{x}^{(n)}) + \frac{1}{N} \sum_m \log(1 - f(\hat{\mathbf{x}}^{(m)})) \end{aligned}$$

Why Jensen-Shannon Divergence? I

- Given a fixed g , we have

$$\begin{aligned} C^* &= \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)}))) \\ &= \max_{\Theta_f} \frac{1}{N} \sum_n \log f(\mathbf{x}^{(n)}) + \frac{1}{N} \sum_m \log(1 - f(\hat{\mathbf{x}}^{(m)})) \\ &\approx \max_{\Theta_f} \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [\log f(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim P_g} [\log(1 - f(\mathbf{x}))] \end{aligned}$$

Why Jensen-Shannon Divergence? I

- Given a fixed g , we have

$$\begin{aligned} C^* &= \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)}))) \\ &= \max_{\Theta_f} \frac{1}{N} \sum_n \log f(\mathbf{x}^{(n)}) + \frac{1}{N} \sum_m \log(1 - f(\hat{\mathbf{x}}^{(m)})) \\ &\approx \max_{\Theta_f} \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [\log f(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim P_g} [\log(1 - f(\mathbf{x}))] \\ &= \max_{\Theta_f} \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{x}} P_g(\mathbf{x}) \log(1 - f(\mathbf{x})) d\mathbf{x} \end{aligned}$$

Why Jensen-Shannon Divergence? I

- Given a fixed g , we have

$$\begin{aligned} C^* &= \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)}))) \\ &= \max_{\Theta_f} \frac{1}{N} \sum_n \log f(\mathbf{x}^{(n)}) + \frac{1}{N} \sum_m \log(1 - f(\hat{\mathbf{x}}^{(m)})) \\ &\approx \max_{\Theta_f} \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [\log f(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim P_g} [\log(1 - f(\mathbf{x}))] \\ &= \max_{\Theta_f} \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{x}} P_g(\mathbf{x}) \log(1 - f(\mathbf{x})) d\mathbf{x} \\ &= \max_{\Theta_f} \int_{\mathbf{x}} [P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) + P_g(\mathbf{x}) \log(1 - f(\mathbf{x}))] d\mathbf{x} \end{aligned}$$

Why Jensen-Shannon Divergence? I

- Given a fixed g , we have

$$\begin{aligned} C^* &= \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log(1 - f(g(\mathbf{c}^{(m)}))) \\ &= \max_{\Theta_f} \frac{1}{N} \sum_n \log f(\mathbf{x}^{(n)}) + \frac{1}{N} \sum_m \log(1 - f(\hat{\mathbf{x}}^{(m)})) \\ &\approx \max_{\Theta_f} \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [\log f(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim P_g} [\log(1 - f(\mathbf{x}))] \\ &= \max_{\Theta_f} \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{x}} P_g(\mathbf{x}) \log(1 - f(\mathbf{x})) d\mathbf{x} \\ &= \max_{\Theta_f} \int_{\mathbf{x}} [P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) + P_g(\mathbf{x}) \log(1 - f(\mathbf{x}))] d\mathbf{x} \end{aligned}$$

- To have C^* , we can find f maximizing

$$P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) + P_g(\mathbf{x}) \log(1 - f(\mathbf{x}))$$

for each \mathbf{x}

- Assuming that f has infinite capacity

Why Jensen-Shannon Divergence? II

Why Jensen-Shannon Divergence? II

- Given P_{data} , P_g , and \mathbf{x} , what is the $f(\mathbf{x})$ that maximizes

$$P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) + P_g(\mathbf{x}) \log(1 - f(\mathbf{x}))?$$

- $f^*(\mathbf{x}) = \frac{P_{\text{data}}(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x})} \in [0, 1]$ [Proof]
- That is,

$$C^* = \max_{\Theta_f} \int_{\mathbf{x}} [P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) + P_g(\mathbf{x}) \log(1 - f(\mathbf{x}))] d\mathbf{x}$$

Why Jensen-Shannon Divergence? II

- Given P_{data} , P_g , and \mathbf{x} , what is the $f(\mathbf{x})$ that maximizes

$$P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) + P_g(\mathbf{x}) \log(1 - f(\mathbf{x}))?$$

- $f^*(\mathbf{x}) = \frac{P_{\text{data}}(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x})} \in [0, 1]$ [Proof]
- That is,

$$\begin{aligned} C^* &= \max_{\Theta_f} \int_{\mathbf{x}} [P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) + P_g(\mathbf{x}) \log(1 - f(\mathbf{x}))] d\mathbf{x} \\ &= \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log \frac{P_{\text{data}}(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x})} d\mathbf{x} \\ &\quad + \int_{\mathbf{x}} P_g(\mathbf{x}) \log \left(\frac{P_g(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x})} \right) d\mathbf{x} \end{aligned}$$

Why Jensen-Shannon Divergence? II

- Given P_{data} , P_g , and \mathbf{x} , what is the $f(\mathbf{x})$ that maximizes

$$P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) + P_g(\mathbf{x}) \log(1 - f(\mathbf{x}))?$$

- $f^*(\mathbf{x}) = \frac{P_{\text{data}}(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x})} \in [0, 1]$ [Proof]
- That is,

$$\begin{aligned} C^* &= \max_{\Theta_f} \int_{\mathbf{x}} [P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) + P_g(\mathbf{x}) \log(1 - f(\mathbf{x}))] d\mathbf{x} \\ &= \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log \frac{P_{\text{data}}(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x})} d\mathbf{x} \\ &\quad + \int_{\mathbf{x}} P_g(\mathbf{x}) \log \left(\frac{P_g(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x})} \right) d\mathbf{x} \\ &= -2 \log 2 + \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log \frac{P_{\text{data}}(\mathbf{x})}{(P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x}))/2} d\mathbf{x} \\ &\quad + \int_{\mathbf{x}} P_g(\mathbf{x}) \log \left(\frac{P_g(\mathbf{x})}{(P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x}))/2} \right) d\mathbf{x} \end{aligned}$$

Why Jensen-Shannon Divergence? II

- Given P_{data} , P_g , and \mathbf{x} , what is the $f(\mathbf{x})$ that maximizes

$$P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) + P_g(\mathbf{x}) \log(1 - f(\mathbf{x}))?$$

- $f^*(\mathbf{x}) = \frac{P_{\text{data}}(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x})} \in [0, 1]$ [Proof]
- That is,

$$\begin{aligned} C^* &= \max_{\Theta_f} \int_{\mathbf{x}} [P_{\text{data}}(\mathbf{x}) \log f(\mathbf{x}) + P_g(\mathbf{x}) \log(1 - f(\mathbf{x}))] d\mathbf{x} \\ &= \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log \frac{P_{\text{data}}(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x})} d\mathbf{x} \\ &\quad + \int_{\mathbf{x}} P_g(\mathbf{x}) \log \left(\frac{P_g(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x})} \right) d\mathbf{x} \\ &= -2 \log 2 + \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log \frac{P_{\text{data}}(\mathbf{x})}{(P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x}))/2} d\mathbf{x} \\ &\quad + \int_{\mathbf{x}} P_g(\mathbf{x}) \log \left(\frac{P_g(\mathbf{x})}{(P_{\text{data}}(\mathbf{x}) + P_g(\mathbf{x}))/2} \right) d\mathbf{x} \\ &= -2 \log 2 + 2D_{JS}(P_{\text{data}} \| P_g) \end{aligned}$$

Balance between g and f , Revisited I

- Cost function of GAN:

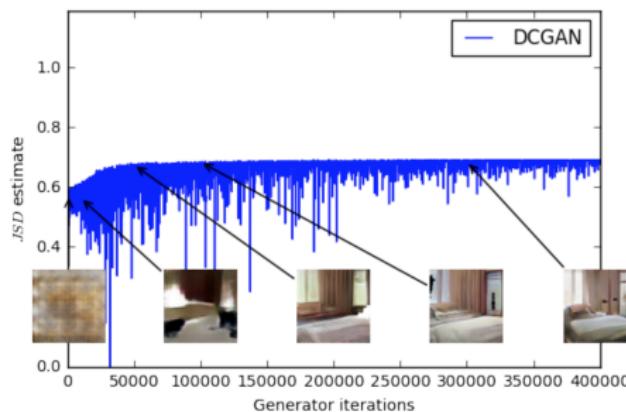
$$\begin{aligned} & \arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log (1 - f(g(\mathbf{c}^{(m)}))) \\ &= \arg \min_{\Theta_g} -2 \log 2 + 2 D_{JS}(P_{\text{data}} \| P_g) \end{aligned}$$

Balance between g and f , Revisited I

- Cost function of GAN:

$$\begin{aligned} & \arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log (1 - f(g(\mathbf{c}^{(m)}))) \\ &= \arg \min_{\Theta_g} -2 \log 2 + 2 D_{JS}(P_{\text{data}} \| P_g) \end{aligned}$$

- However, no matter how g changes, $D_{JS}(P_{\text{data}} \| P_g)$ remains high during the GAN training process

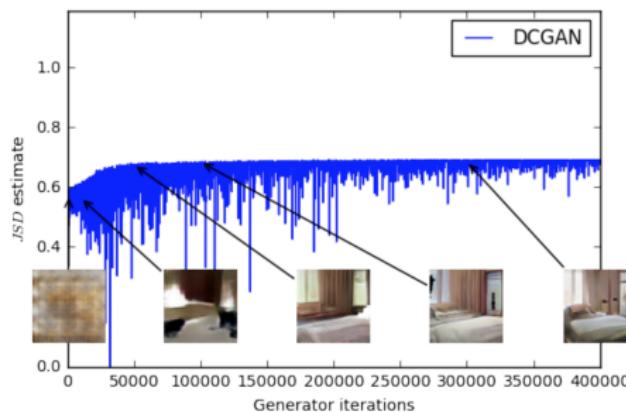


Balance between g and f , Revisited I

- Cost function of GAN:

$$\begin{aligned} & \arg \min_{\Theta_g} \max_{\Theta_f} \sum_n \log f(\mathbf{x}^{(n)}) + \sum_m \log (1 - f(g(\mathbf{c}^{(m)}))) \\ &= \arg \min_{\Theta_g} -2 \log 2 + 2 D_{JS}(P_{\text{data}} \| P_g) \end{aligned}$$

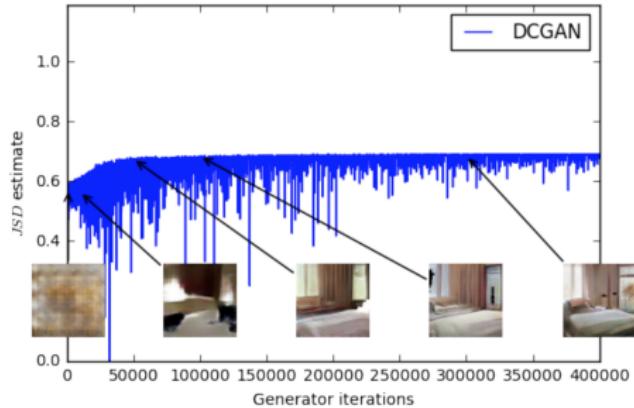
- However, no matter how g changes, $D_{JS}(P_{\text{data}} \| P_g)$ remains high during the GAN training process



- There's something wrong with the design of the inner max problem!

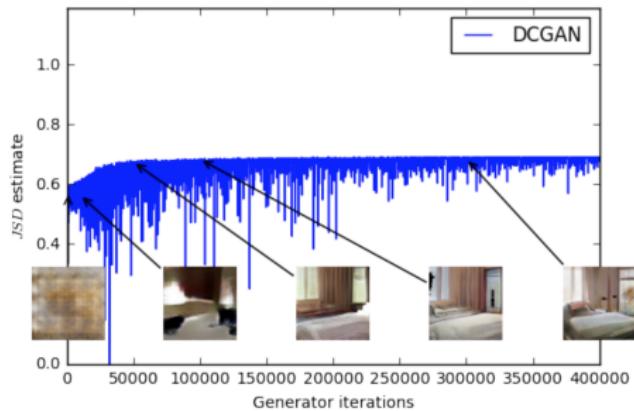
Balance between g and f , Revisited II

- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$



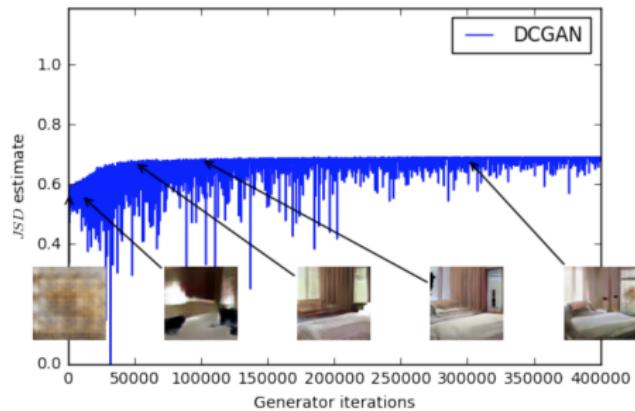
Balance between g and f , Revisited II

- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$
- $0 \leq D_{JS}(P_{\text{data}} \| P_g) \leq \log 2 \approx 0.69$



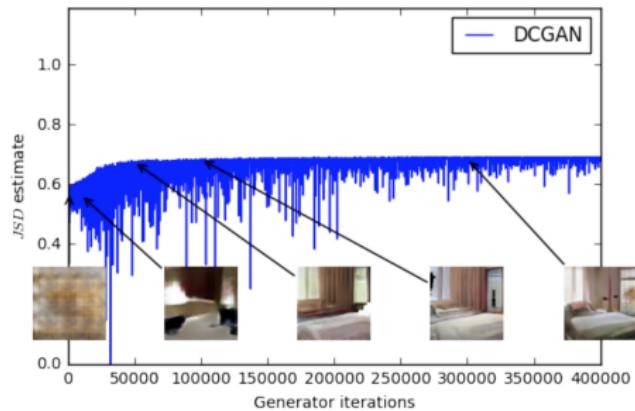
Balance between g and f , Revisited II

- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$
- $0 \leq D_{JS}(P_{\text{data}} \| P_g) \leq \log 2 \approx 0.69$
- When does $D_{JS}(P_{\text{data}} \| P_g)$ reach its maximum value?



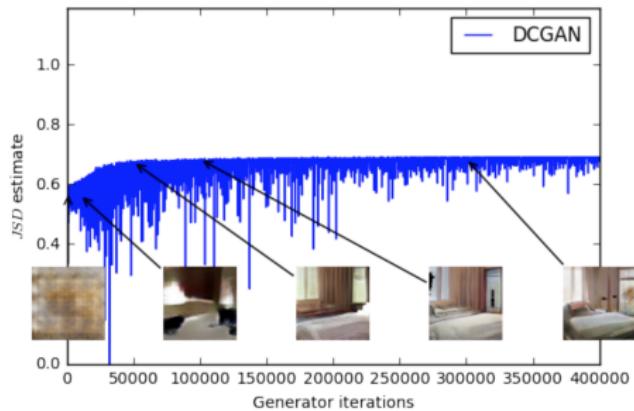
Balance between g and f , Revisited II

- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$
- $0 \leq D_{JS}(P_{\text{data}} \| P_g) \leq \log 2 \approx 0.69$
- When does $D_{JS}(P_{\text{data}} \| P_g)$ reach its maximum value?
- Let's see the case when P_g and P_{data} are “disjointed”



Balance between g and f , Revisited II

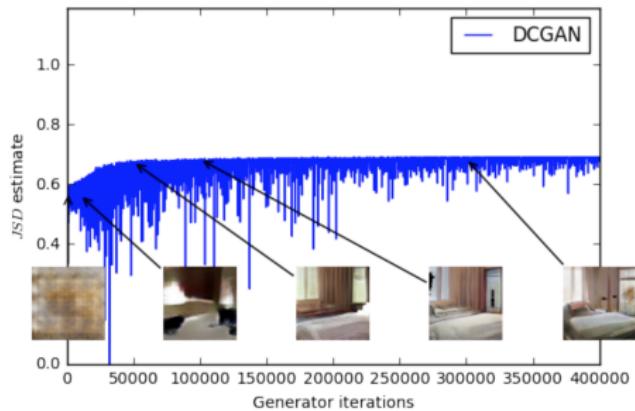
- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$
- $0 \leq D_{JS}(P_{\text{data}} \| P_g) \leq \log 2 \approx 0.69$
- When does $D_{JS}(P_{\text{data}} \| P_g)$ reach its maximum value?
- Let's see the case when P_g and P_{data} are “disjointed”
- Suppose $P_g(x) \neq 0 \Leftrightarrow P_{\text{data}}(x) = 0$ and $P_g(x) = 0 \Leftrightarrow P_{\text{data}}(x) \neq 0$, we have



$$D_{JS}(P_g \| P_{\text{data}}) = \frac{1}{2} D_{KL}(P_g \| \frac{P_g + P_{\text{data}}}{2}) + \frac{1}{2} D_{KL}(P_{\text{data}} \| \frac{P_g + P_{\text{data}}}{2})$$

Balance between g and f , Revisited II

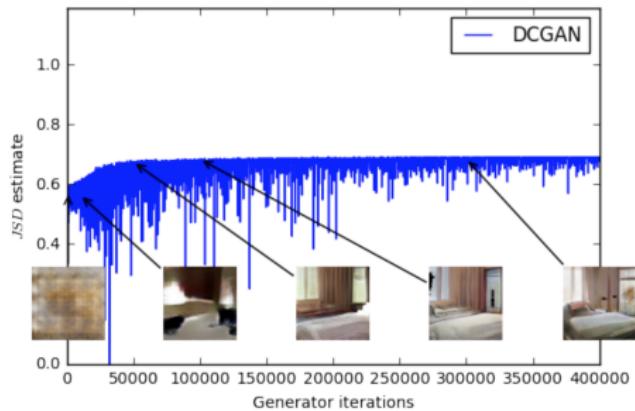
- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$
- $0 \leq D_{JS}(P_{\text{data}} \| P_g) \leq \log 2 \approx 0.69$
- When does $D_{JS}(P_{\text{data}} \| P_g)$ reach its maximum value?
- Let's see the case when P_g and P_{data} are “disjointed”
- Suppose $P_g(\mathbf{x}) \neq 0 \Leftrightarrow P_{\text{data}}(\mathbf{x}) = 0$ and $P_g(\mathbf{x}) = 0 \Leftrightarrow P_{\text{data}}(\mathbf{x}) \neq 0$, we have



$$\begin{aligned} D_{JS}(P_g \| P_{\text{data}}) &= \frac{1}{2} D_{KL}(P_g \| \frac{P_g + P_{\text{data}}}{2}) + \frac{1}{2} D_{KL}(P_{\text{data}} \| \frac{P_g + P_{\text{data}}}{2}) \\ &= \frac{1}{2} \int_{\mathbf{x}} P_g(\mathbf{x}) \log \frac{2P_g(\mathbf{x})}{P_g(\mathbf{x}) + P_{\text{data}}(\mathbf{x})} d\mathbf{x} \\ &\quad + \frac{1}{2} \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log \frac{2P_{\text{data}}(\mathbf{x})}{P_g(\mathbf{x}) + P_{\text{data}}(\mathbf{x})} d\mathbf{x} \end{aligned}$$

Balance between g and f , Revisited II

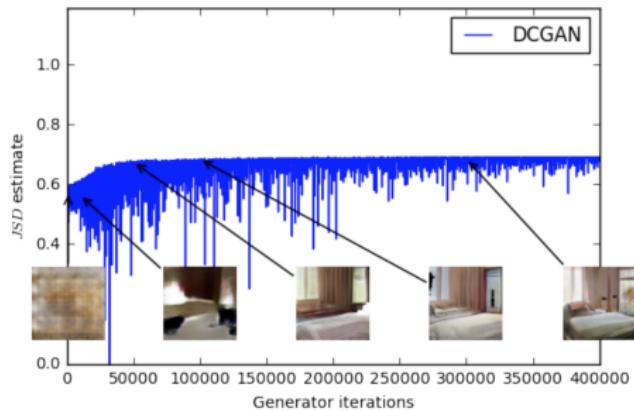
- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$
- $0 \leq D_{JS}(P_{\text{data}} \| P_g) \leq \log 2 \approx 0.69$
- When does $D_{JS}(P_{\text{data}} \| P_g)$ reach its maximum value?
- Let's see the case when P_g and P_{data} are “disjointed”
- Suppose $P_g(\mathbf{x}) \neq 0 \Leftrightarrow P_{\text{data}}(\mathbf{x}) = 0$ and $P_g(\mathbf{x}) = 0 \Leftrightarrow P_{\text{data}}(\mathbf{x}) \neq 0$, we have



$$\begin{aligned} D_{JS}(P_g \| P_{\text{data}}) &= \frac{1}{2} D_{KL}(P_g \| \frac{P_g + P_{\text{data}}}{2}) + \frac{1}{2} D_{KL}(P_{\text{data}} \| \frac{P_g + P_{\text{data}}}{2}) \\ &= \frac{1}{2} \int_{\mathbf{x}} P_g(\mathbf{x}) \log \frac{2P_g(\mathbf{x})}{P_g(\mathbf{x}) + P_{\text{data}}(\mathbf{x})} d\mathbf{x} \\ &\quad + \frac{1}{2} \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log \frac{2P_{\text{data}}(\mathbf{x})}{P_g(\mathbf{x}) + P_{\text{data}}(\mathbf{x})} d\mathbf{x} \\ &= \frac{1}{2} \log 2 + \frac{1}{2} \log 2 = \log 2 \end{aligned}$$

Balance between g and f , Revisited II

- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$
- $0 \leq D_{JS}(P_{\text{data}} \| P_g) \leq \log 2 \approx 0.69$
- When does $D_{JS}(P_{\text{data}} \| P_g)$ reach its maximum value?
- Let's see the case when P_g and P_{data} are “disjointed”
- Suppose $P_g(x) \neq 0 \Leftrightarrow P_{\text{data}}(x) = 0$ and $P_g(x) = 0 \Leftrightarrow P_{\text{data}}(x) \neq 0$, we have

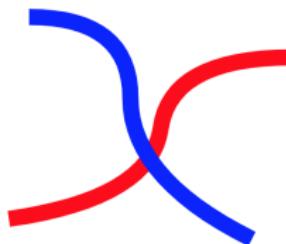


$$\begin{aligned} D_{JS}(P_g \| P_{\text{data}}) &= \frac{1}{2} D_{KL}(P_g \| \frac{P_g + P_{\text{data}}}{2}) + \frac{1}{2} D_{KL}(P_{\text{data}} \| \frac{P_g + P_{\text{data}}}{2}) \\ &= \frac{1}{2} \int_x P_g(x) \log \frac{2P_g(x)}{P_g(x) + P_{\text{data}}(x)} dx \\ &\quad + \frac{1}{2} \int_x P_{\text{data}}(x) \log \frac{2P_{\text{data}}(x)}{P_g(x) + P_{\text{data}}(x)} dx \\ &= \frac{1}{2} \log 2 + \frac{1}{2} \log 2 = \log 2 \end{aligned}$$

- Are P_g and P_{data} really disjointed during the GAN training?

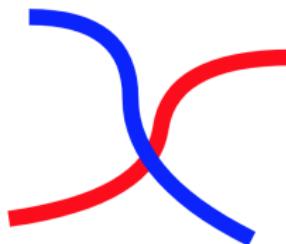
Disjoining P_g and P_{data}

- In a high dimensional space, \mathbf{x} and $g(\mathbf{z})$ may reside in low dimensional manifolds
 - P_g and P_{data} may have values only on the manifolds



Disjoining P_g and P_{data}

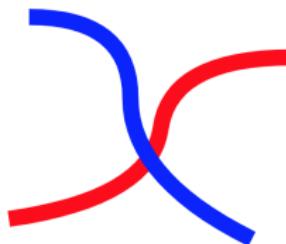
- In a high dimensional space, \mathbf{x} and $g(\mathbf{z})$ may reside in low dimensional manifolds
 - P_g and P_{data} may have values only on the manifolds



- P_g and P_{data} can be very different initially during GAN training

Disjoining P_g and P_{data}

- In a high dimensional space, \mathbf{x} and $g(\mathbf{z})$ may reside in low dimensional manifolds
 - P_g and P_{data} may have values only on the manifolds



- P_g and P_{data} can be very different initially during GAN training
- The intersections where $P_g(\mathbf{x}) \neq 0$ and $P_{\text{data}}(\mathbf{x}) \neq 0$ can be neglected
 - $P_g(\mathbf{x}) \neq 0 \Leftrightarrow P_{\text{data}}(\mathbf{x}) = 0$ and $P_g(\mathbf{x}) = 0 \Leftrightarrow P_{\text{data}}(\mathbf{x}) \neq 0$ almost surely
 - Maximum JS divergence at all time

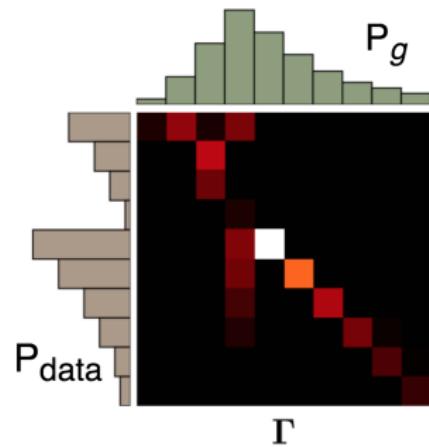
Better Divergence Measure?

- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$
 - $D_{JS}(P_{\text{data}} \| P_g)$ does **not** bring P_g and P_{data} closer during GAN training

Better Divergence Measure?

- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$
 - $D_{JS}(P_{\text{data}} \| P_g)$ does **not** bring P_g and P_{data} closer during GAN training
- Wasserstein (or earth-mover) distance:

$$\begin{aligned} W(P_{\text{data}}, P_g) &= \inf_{Q \in \Gamma(P_{\text{data}}, P_g)} E_{(x, \hat{x}) \sim Q} [\|x - \hat{x}\|] \\ &= \inf_{Q \in \Gamma(P_{\text{data}}, P_g)} \int_{(x, \hat{x})} Q(x, \hat{x}) \|x - \hat{x}\| d(x, \hat{x}) \end{aligned}$$

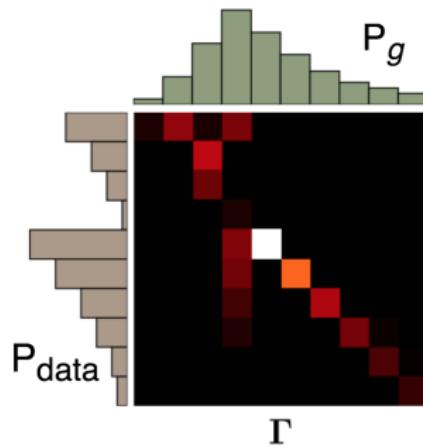


Better Divergence Measure?

- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$
 - $D_{JS}(P_{\text{data}} \| P_g)$ does **not** bring P_g and P_{data} closer during GAN training
- Wasserstein (or earth-mover) distance:

$$\begin{aligned} W(P_{\text{data}}, P_g) &= \inf_{Q \in \Gamma(P_{\text{data}}, P_g)} E_{(x, \hat{x}) \sim Q} [\|x - \hat{x}\|] \\ &= \inf_{Q \in \Gamma(P_{\text{data}}, P_g)} \int_{(x, \hat{x})} Q(x, \hat{x}) \|x - \hat{x}\| d(x, \hat{x}) \end{aligned}$$

- Intuitively, the minimal “cost” to change P_{data} into P_g

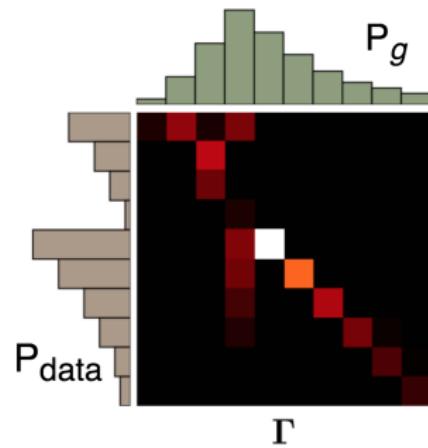


Better Divergence Measure?

- GAN: $\arg \min_{\Theta_g} D_{JS}(P_{\text{data}} \| P_g)$
 - $D_{JS}(P_{\text{data}} \| P_g)$ does **not** bring P_g and P_{data} closer during GAN training
- Wasserstein (or earth-mover) distance:

$$\begin{aligned} W(P_{\text{data}}, P_g) &= \inf_{Q \in \Gamma(P_{\text{data}}, P_g)} E_{(x, \hat{x}) \sim Q} [\|x - \hat{x}\|] \\ &= \inf_{Q \in \Gamma(P_{\text{data}}, P_g)} \int_{(x, \hat{x})} Q(x, \hat{x}) \|x - \hat{x}\| d(x, \hat{x}) \end{aligned}$$

- Intuitively, the minimal “cost” to change P_{data} into P_g
- $W(P_{\text{data}}, P_g)$ measures the “divergence” between P_g and P_{data} **even when they are disjointed**



Wasserstein GAN I

- W-GAN: $\arg \min_{\Theta_g} W(P_{\text{data}}, P_g)$
 - $W(P_{\text{data}}, P_g) = \inf_{Q \in \Gamma(P_{\text{data}}, P_g)} E_{(\mathbf{x}, \hat{\mathbf{x}}) \sim Q} [\|\mathbf{x} - \hat{\mathbf{x}}\|]$

Wasserstein GAN I

- W-GAN: $\arg \min_{\Theta_g} W(P_{\text{data}}, P_g)$
 - $W(P_{\text{data}}, P_g) = \inf_{Q \in \Gamma(P_{\text{data}}, P_g)} E_{(\mathbf{x}, \hat{\mathbf{x}}) \sim Q} [\|\mathbf{x} - \hat{\mathbf{x}}\|]$
- Unfortunately, $W(P_{\text{data}}, P_g)$ is hard to solve directly

Wasserstein GAN II

Theorem

Consider f 's that are Lipschitz continuous with constant 1, i.e.,

$$|f(\mathbf{x}) - f(\hat{\mathbf{x}})| \leq 1 \cdot \|\mathbf{x} - \hat{\mathbf{x}}\|, \forall \mathbf{x}, \hat{\mathbf{x}},$$

we have ^a

$$W(P_{\text{data}}, P_g) = \sup_f E_{\mathbf{x} \sim P_{\text{data}}} [f(\mathbf{x})] - E_{\mathbf{x} \sim P_g} [f(\mathbf{x})]$$

Wasserstein GAN II

Theorem

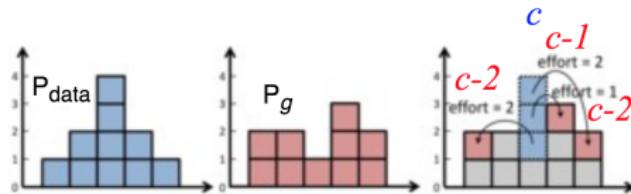
Consider f 's that are Lipschitz continuous with constant 1, i.e.,

$$|f(\mathbf{x}) - f(\hat{\mathbf{x}})| \leq 1 \cdot \|\mathbf{x} - \hat{\mathbf{x}}\|, \forall \mathbf{x}, \hat{\mathbf{x}},$$

we have ^a

$$\begin{aligned} W(P_{\text{data}}, P_g) &= \sup_f E_{\mathbf{x} \sim P_{\text{data}}} [f(\mathbf{x})] - E_{\mathbf{x} \sim P_g} [f(\mathbf{x})] \\ &= \sup_f \int_{\mathbf{x}} (P_{\text{data}}(\mathbf{x}) - P_g(\mathbf{x})) f(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

^a<https://vincentherrmann.github.io/blog/wasserstein/>



Wasserstein GAN II

Theorem

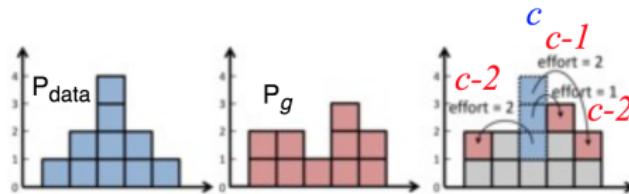
Consider f 's that are Lipschitz continuous with constant 1, i.e.,

$$|f(\mathbf{x}) - f(\hat{\mathbf{x}})| \leq 1 \cdot \|\mathbf{x} - \hat{\mathbf{x}}\|, \forall \mathbf{x}, \hat{\mathbf{x}},$$

we have ^a

$$\begin{aligned} W(P_{\text{data}}, P_g) &= \sup_f E_{\mathbf{x} \sim P_{\text{data}}} [f(\mathbf{x})] - E_{\mathbf{x} \sim P_g} [f(\mathbf{x})] \\ &= \sup_f \int_{\mathbf{x}} (P_{\text{data}}(\mathbf{x}) - P_g(\mathbf{x})) f(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

^a<https://vincentherrmann.github.io/blog/wasserstein/>



- W-GAN [1]: $\arg \min_{\Theta_g} \max_{\Theta_f} E_{\mathbf{x} \sim P_{\text{data}}} [f(\mathbf{x})] - E_{\mathbf{x} \sim P_g} [f(\mathbf{x})]$

Alternate SGD for W-GAN

$$\begin{aligned} & \arg \min_{\Theta_g} \max_{\Theta_f} \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim P_g} [f(\mathbf{x})] \\ &= \arg \min_{\Theta_g} \max_{\Theta_f} \sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)})) \end{aligned}$$

- f a regressor **without** the sigmoid output layer

Alternate SGD for W-GAN

$$\begin{aligned} & \arg \min_{\Theta_g} \max_{\Theta_f} \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim P_g} [f(\mathbf{x})] \\ &= \arg \min_{\Theta_g} \max_{\Theta_f} \sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)})) \end{aligned}$$

- f a regressor **without** the sigmoid output layer
- Initialize Θ_g for g and Θ_f for f
- At each SGD step/iteration:
 - ① Repeat K times (with fixed Θ_g):
 - ① Sample N real points $\{\mathbf{x}^{(n)}\}_n$ from \mathbb{X} and N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - ② $\Theta_f \leftarrow \Theta_f + \eta \text{clip}(\nabla_{\Theta_f} [\sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)}))])$
 - ② Execute once (with fixed Θ_f):
 - ① Sample N codes from $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - ② $\Theta_g \leftarrow \Theta_g - \eta \nabla_{\Theta_g} [-\sum_m f(g(\mathbf{c}^{(m)}))]$

Why Gradient Clipping?

- Update rule for Θ_f :

$$\Theta_f \leftarrow \Theta_f + \eta \text{clip}(\nabla_{\Theta_f} [\sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)}))])$$

- Gradient clipping: $\forall w \in \Theta_f$, $\text{clip}(w) = \max(\min(w, \tau), -\tau)$ for some threshold $\tau > 0$
- Why?

Why Gradient Clipping?

- Update rule for Θ_f :

$$\Theta_f \leftarrow \Theta_f + \eta \text{clip}(\nabla_{\Theta_f} [\sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)}))])$$

- Gradient clipping: $\forall w \in \Theta_f$, $\text{clip}(w) = \max(\min(w, \tau), -\tau)$ for some threshold $\tau > 0$
- Why?
- In W-GAN, we have :

$$\begin{aligned} & \arg \min_{\Theta_g} W(P_{\text{data}}, P_g) \\ &= \arg \min_{\Theta_g} \max_{\Theta_f} \sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)})) \end{aligned}$$

only if f is 1-Lipschitz continuous

Why Gradient Clipping?

- Update rule for Θ_f :

$$\Theta_f \leftarrow \Theta_f + \eta \text{clip}(\nabla_{\Theta_f} [\sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)}))])$$

- Gradient clipping: $\forall w \in \Theta_f$, $\text{clip}(w) = \max(\min(w, \tau), -\tau)$ for some threshold $\tau > 0$
- Why?
- In W-GAN, we have :

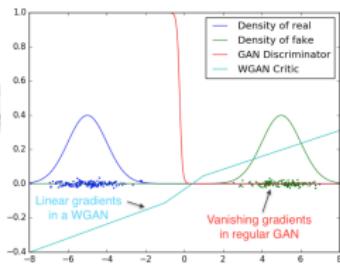
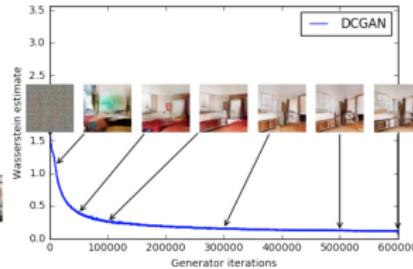
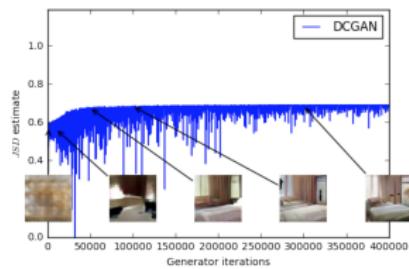
$$\begin{aligned} & \arg \min_{\Theta_g} W(P_{\text{data}}, P_g) \\ &= \arg \min_{\Theta_g} \max_{\Theta_f} \sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)})) \end{aligned}$$

only if f is 1-Lipschitz continuous

- Heuristics for making f 1-Lipchitz

Advantages of W-GAN

- “Corrects” the inner max problem
 - Wasserstein distance guides g even when P_g and P_{data} “disjointed”
 - Training less sensitive to K (balance between g and f)
- The max value can be used as a “stop” indicator
- f a regressor, avoids vanishing gradients for g



Improved W-GAN I

- In practice, W-GAN training converges slowly and is unstable to τ

Improved W-GAN I

- In practice, W-GAN training converges slowly and is unstable to τ
- W-GAN use a small τ to make f 1-Lipschitz continuous
- However, too small a τ severely limits the capacity of f such it cannot actually maximize

$$\max_{\Theta_f} \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim P_g} [f(\mathbf{x})]$$

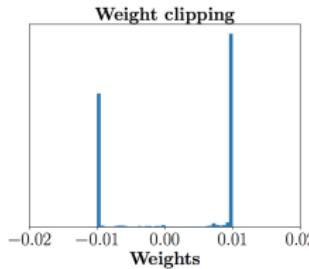
- g is not updated for minimizing $W(P_{\text{data}}, P_g)$

Improved W-GAN I

- In practice, W-GAN training converges slowly and is unstable to τ
- W-GAN use a small τ to make f 1-Lipschitz continuous
- However, too small a τ severely limits the capacity of f such it cannot actually maximize

$$\max_{\Theta_f} \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim P_g} [f(\mathbf{x})]$$

- g is not updated for minimizing $W(P_{\text{data}}, P_g)$
- Distribution of weight values of f ($\tau = 0.01$):



- Exploding and vanishing gradients

Improved W-GAN II

- If f is 1-Lipschitz, then $\|\nabla f(\mathbf{x})\| \leq 1$ for all \mathbf{x}
- Why not just penalize $\|\nabla f(\mathbf{x})\| > 1$ for all \mathbf{x} ?
- Cost function:

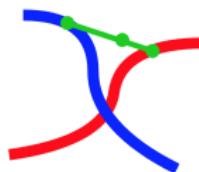
$$\begin{aligned} \arg \min_{\Theta_g} \max_{\Theta_f} E_{\mathbf{x} \sim P_{\text{data}}} [f(\mathbf{x})] - E_{\mathbf{x} \sim P_g} [f(\mathbf{x})] \\ - \lambda E_{\mathbf{x} \sim P_{\text{penalty}}} [\max(0, \|\nabla f(\mathbf{x})\| - 1)] \end{aligned}$$

Improved W-GAN II

- If f is 1-Lipschitz, then $\|\nabla f(\mathbf{x})\| \leq 1$ for all \mathbf{x}
- Why not just penalize $\|\nabla f(\mathbf{x})\| > 1$ for all \mathbf{x} ?
- Cost function:

$$\arg \min_{\Theta_g} \max_{\Theta_f} \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim P_g} [f(\mathbf{x})] \\ - \lambda \mathbb{E}_{\mathbf{x} \sim P_{\text{penalty}}} [\max(0, \|\nabla f(\mathbf{x})\| - 1)]$$

- P_{penalty} ?

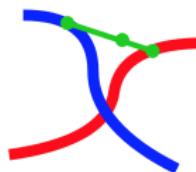


Improved W-GAN II

- If f is 1-Lipschitz, then $\|\nabla f(\mathbf{x})\| \leq 1$ for all \mathbf{x}
- Why not just penalize $\|\nabla f(\mathbf{x})\| > 1$ for all \mathbf{x} ?
- Cost function:

$$\begin{aligned} \arg \min_{\Theta_g} \max_{\Theta_f} \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim P_g} [f(\mathbf{x})] \\ - \lambda \mathbb{E}_{\mathbf{x} \sim P_{\text{penalty}}} [\max(0, \|\nabla f(\mathbf{x})\| - 1)] \end{aligned}$$

- P_{penalty} ?



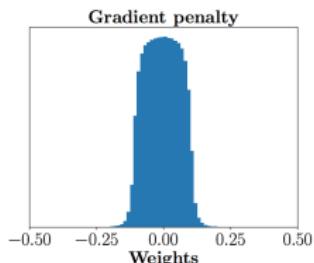
- W-GAN-GP [9]:

$$\arg \min_{\Theta_g} \max_{\Theta_f} \sum_n f(\mathbf{x}^{(n)}) - \sum_m f(g(\mathbf{c}^{(m)})) - \lambda \sum_p (\|\nabla f(\mathbf{x}^{(p)})\| - 1)^2$$

- The larger $f(\mathbf{x}^{(p)})$ the better (subject to $\|\nabla f(\mathbf{x}^{(p)})\| \leq 1$)
- Faster convergence

Results

DCGAN	LSGAN	WGAN (clipping)	WGAN-GP (ours)
Baseline (G : DCGAN, D : DCGAN)			
G : No BN and a constant number of filters, D : DCGAN	 		
G : 4-layer 512-dim ReLU MLP, D : DCGAN	 		
No normalization in either G or D	 		
Gated multiplicative nonlinearities everywhere in G and D	 		
$tanh$ nonlinearities everywhere in G and D	 		
101-layer ResNet G and D	 		



Challenge: Global Coherence

Challenge: Global Coherence

- Large images generated by GANs usually lack global coherency

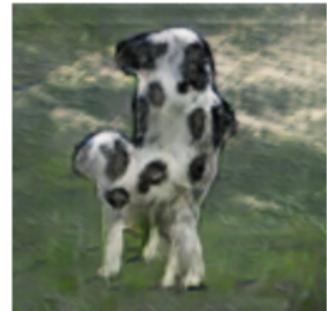
Counting



Perspective



Shape



Challenge: Global Coherence

- Large images generated by GANs usually lack global coherency

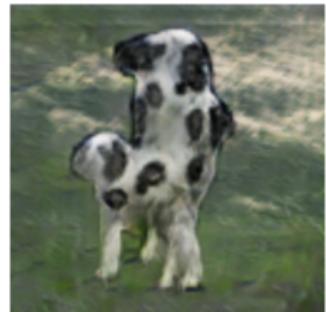
Counting



Perspective



Shape



- A CNN, when used as f , detects *existence* of patterns more than their *relative positions*
 - f loses track of the position of a pattern after several pooling layers
 - Relative position of patterns in \mathbb{X} may change due to different view angels
- Solutions?

Challenge: Global Coherence

- Large images generated by GANs usually lack global coherency

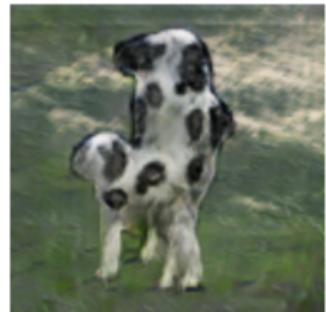
Counting



Perspective

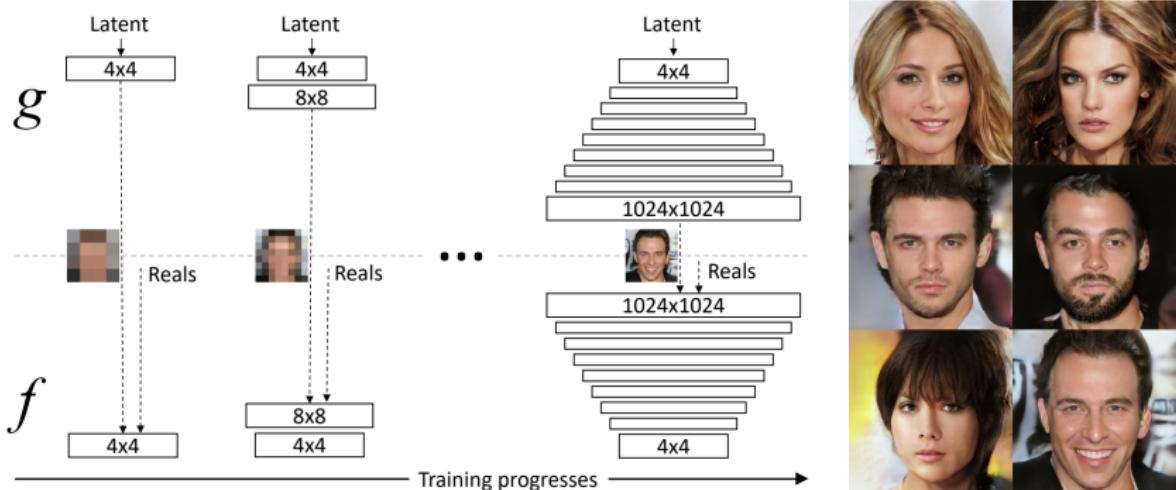


Shape



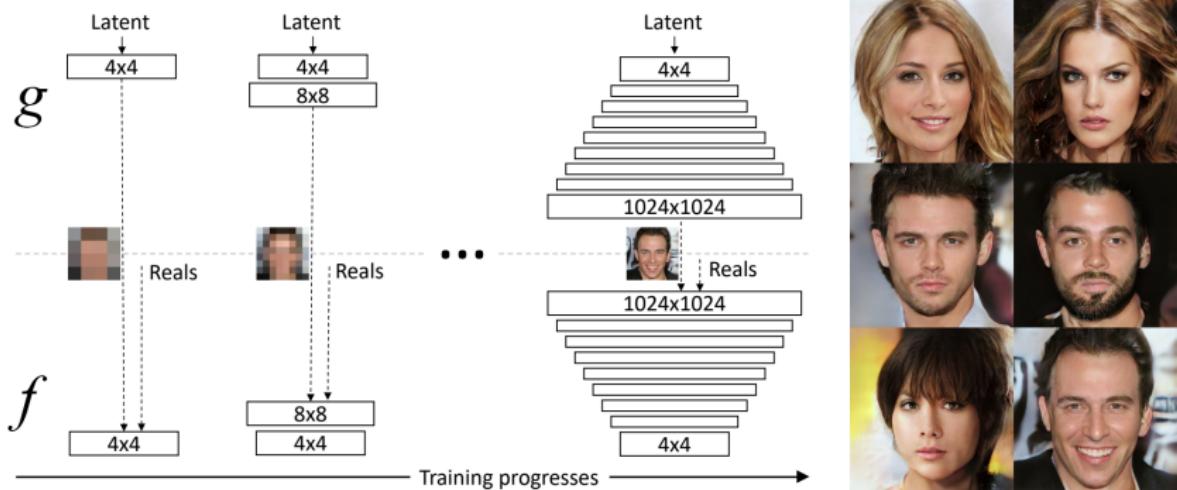
- A CNN, when used as f , detects *existence* of patterns more than their *relative positions*
 - f loses track of the position of a pattern after several pooling layers
 - Relative position of patterns in \mathbb{X} may change due to different view angels
- Solutions? A better f , such as the CapsuleNet [34]?

Progressive Growing of GANs [14]



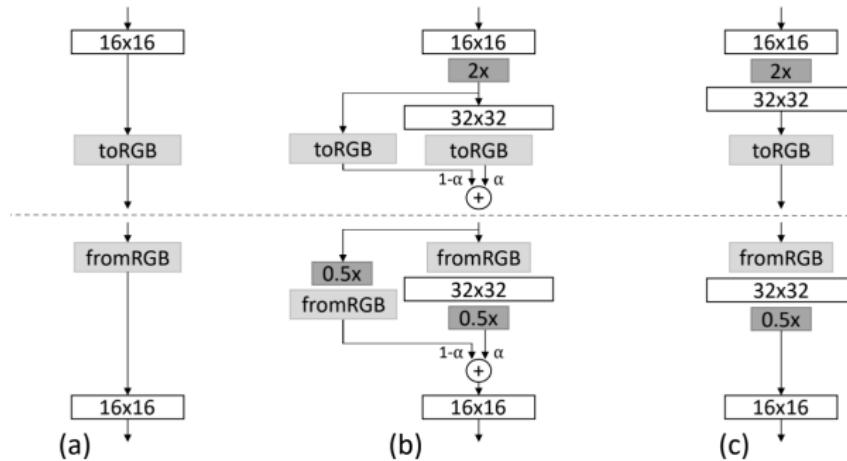
- Incrementally adds new layers in sequential GAN trainings
 - Convolution + upsampling for g each time
 - Convolution + downpooling for f each time
- Real images are downscaled to match the current resolution of g

Progressive Growing of GANs [14]



- Incrementally adds new layers in sequential GAN trainings
 - Convolution + upsampling for g each time
 - Convolution + downpooling for f each time
- Real images are downscaled to match the current resolution of g
- Goal: to let new layers ***add details without ruining the context***

Transition when Adding a New Layer



- Gradually increases α
 - New convolution layer in g/f learns to generate/detect details first
 - Then learns the “context”

Results



- Minibatch discrimination [35] + W-GAN-GP [9] + progressive growing [14] + other tricks
- 2 ~ 6 times faster

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

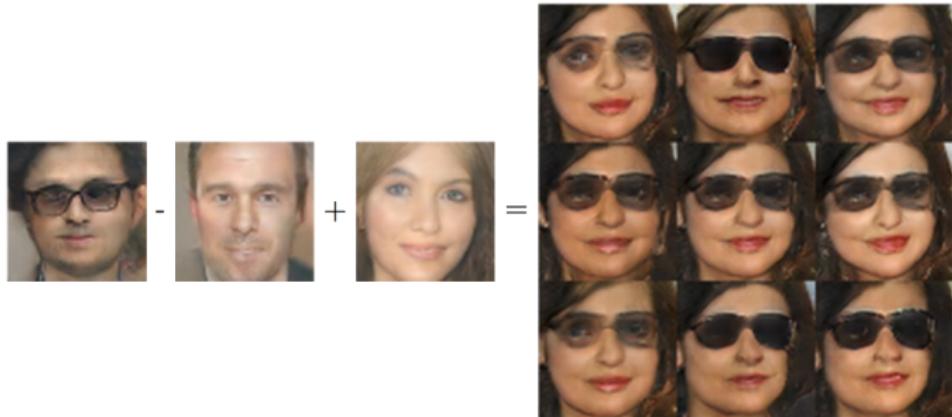
⑥ Diffusion Models

⑦ Generative Adversarial Networks*

- Basic Architecture
- Challenges
- More GANs

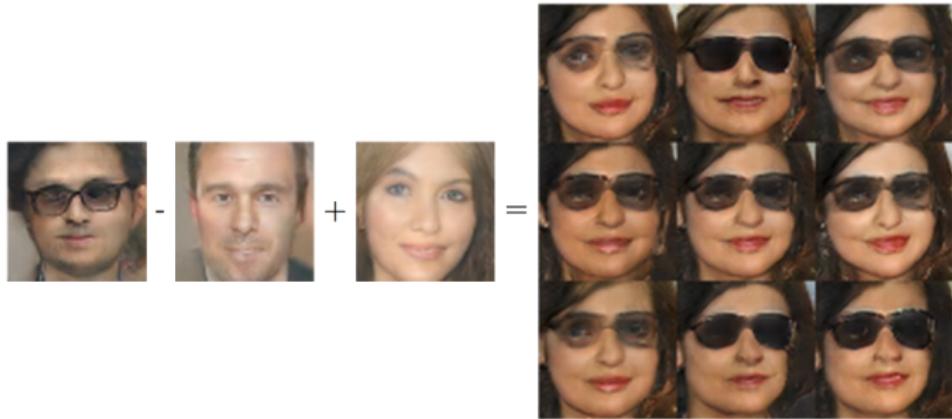
Code Space Arithmetics

- DC-GAN [29] can learn to use codes in meaningful ways:



Code Space Arithmetics

- DC-GAN [29] can learn to use codes in meaningful ways:



- Finding codes for images with constraints [44, 3] [▶ Demo 1](#) [▶ Demo 2](#)

$$\arg \min_{\mathbf{c}} \|mask(g(\mathbf{c})) - \text{constraint}\|_F$$

Conditional GAN I

- Text to image synthesis [32]: $\mathbb{X} = \{(\mathbf{x}^{(n)}, \phi^{(n)})\}_n$
"This bird is completely red with black wings and pointy beak."

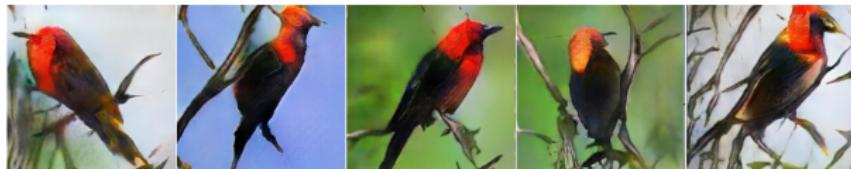


- How?

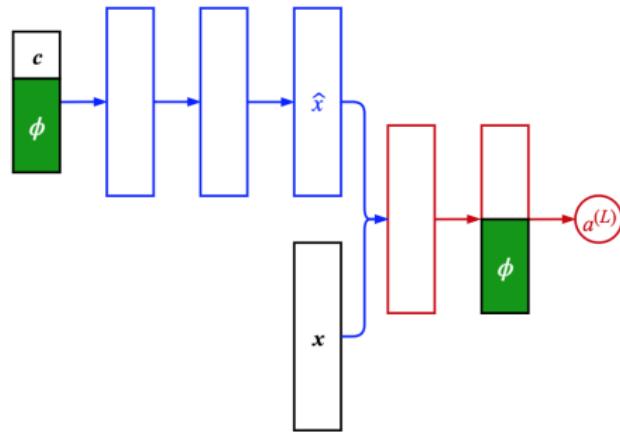
Conditional GAN I

- Text to image synthesis [32]: $\mathbb{X} = \{(x^{(n)}, \phi^{(n)})\}_n$

“This bird is completely red with black wings and pointy beak.”

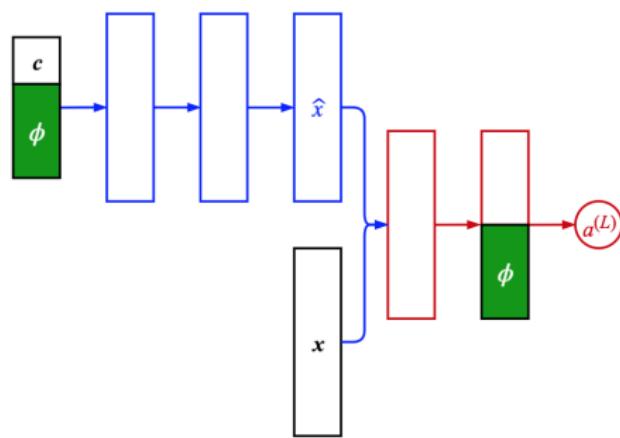


- How?



Conditional GAN II

- Pitfall: g and f can choose to ignore the condition ϕ altogether
- Solution?



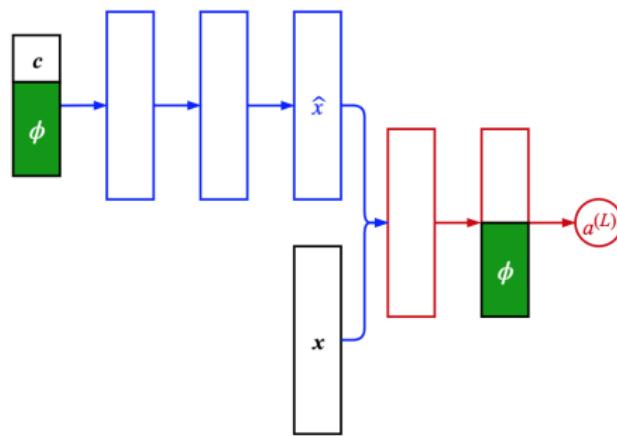
Conditional GAN II

- Pitfall: g and f can choose to ignore the condition ϕ altogether
- Solution? Conditioned labeling

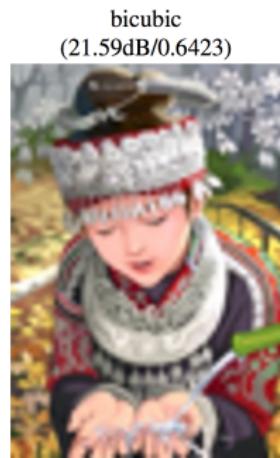
$$(\mathbf{x}^{(n)}, \phi^{(n)}) \Rightarrow \text{true}$$

$$(\mathbf{x}^{(n)}, \phi') \Rightarrow \text{false}, \forall \phi' \neq \phi^{(n)}$$

$$(\hat{\mathbf{x}}^{(m)}, \phi^{(m)}) \Rightarrow \text{false}$$

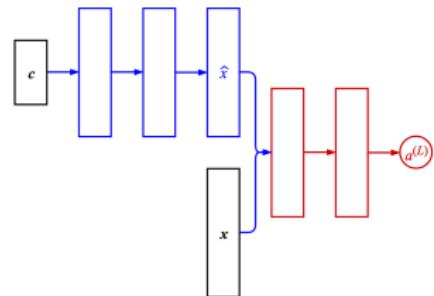


Super Resolution [18]



Super Resolution [18]

- g : low res img \rightarrow high res img
 - Training: c 's are downscaled images



Super Resolution [18]

- g : low res img \rightarrow high res img
 - Training: c 's are downscaled images
- No “creativity,” f acts as a better loss metric

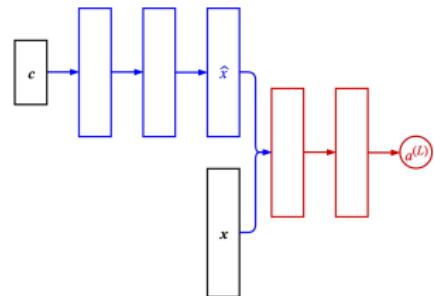


Image-to-Image Translation [13] I

- Given an image x_{src} in source domain, generate image(s) x_{target} in target domain
 - x_{src} and x_{target} are semantically aligned

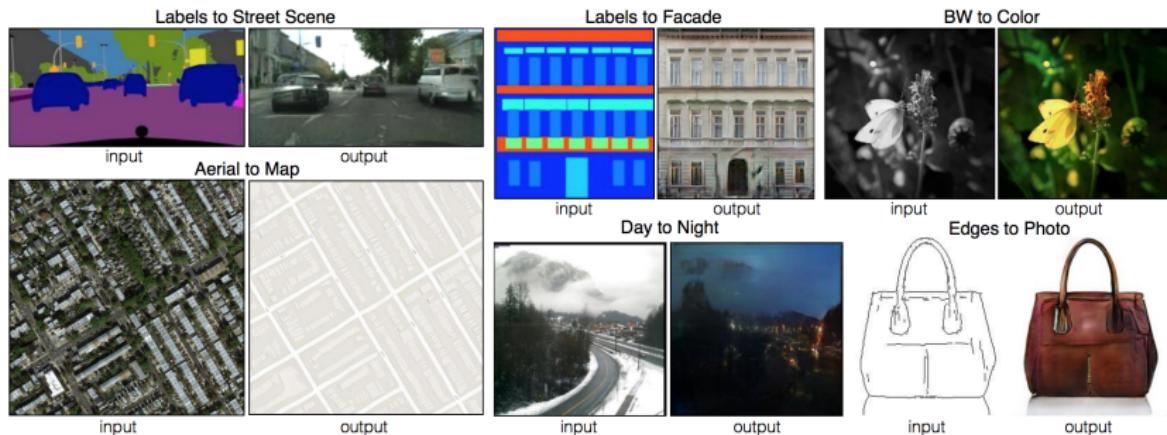


Image-to-Image Translation [13] II

- Based on conditional GAN:
 - $\mathbf{c} = \mathbf{x}_{\text{src}}$; $g(\mathbf{c}) = \mathbf{x}_{\text{target}}$
 - Conditioned labels
 - Uses dropout layers to create diversity (if needed)

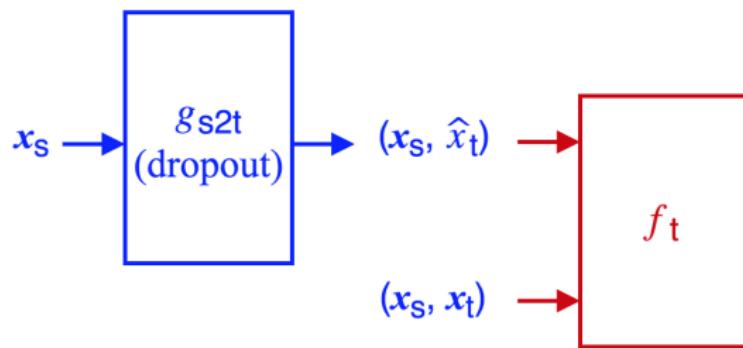
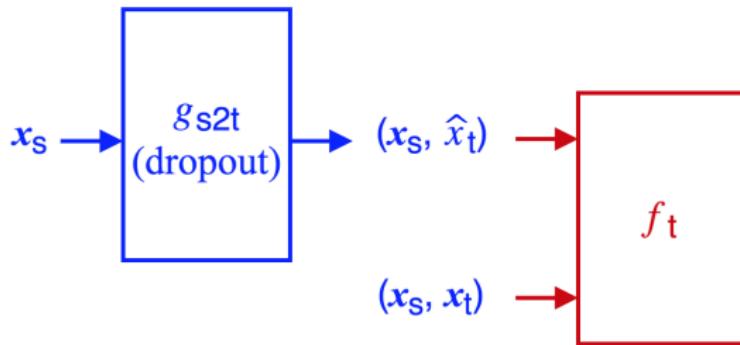


Image-to-Image Translation [13] II

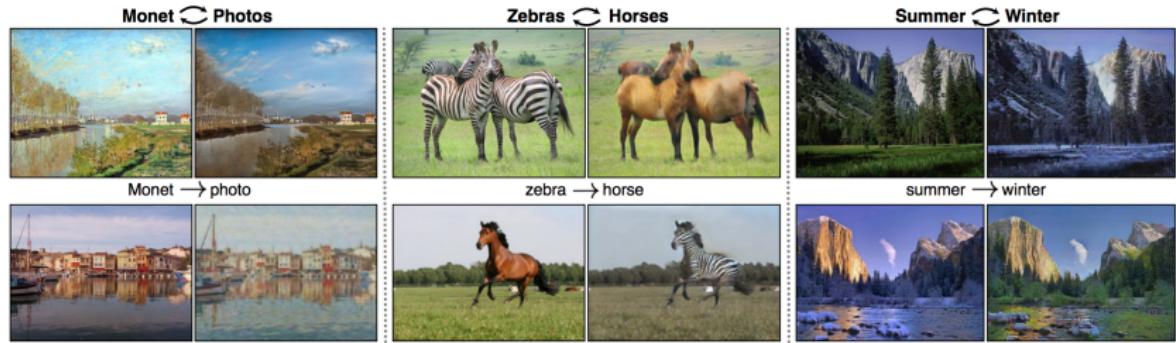
- Based on conditional GAN:
 - $\mathbf{c} = \mathbf{x}_{\text{src}}$; $g(\mathbf{c}) = \mathbf{x}_{\text{target}}$
 - Conditioned labels
 - Uses dropout layers to create diversity (if needed)
- Requires **paired** examples $\mathbb{X} = \{(\mathbf{x}_{\text{src}}^{(n)}, \mathbf{x}_{\text{target}}^{(n)})\}_n$



Unpaired Image-to-Image Translation I

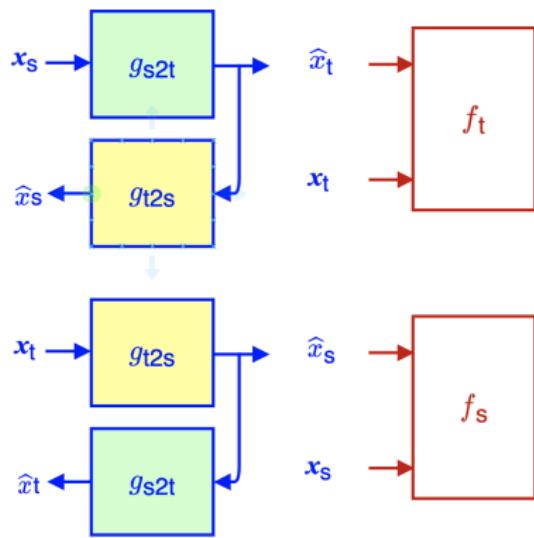
- What if the images in different domains are *unpaired*?

- $\mathbb{X} = \{\mathbf{x}_{\text{src}}^{(n)}\}_n \cup \{\mathbf{x}_{\text{target}}^{(n)}\}_n$



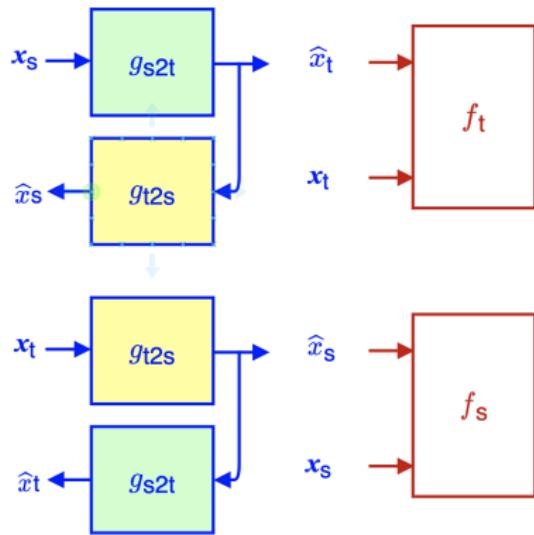
Unpaired Image-to-Image Translation II

- Cycle GAN [45]: to train two generators $g_{\text{src2target}}$ and $g_{\text{target2src}}$ simultaneously in two GANs



Unpaired Image-to-Image Translation II

- Cycle GAN [45]: to train two generators $g_{\text{src2target}}$ and $g_{\text{target2src}}$ simultaneously in two GANs
- Add a loss term $\sum_n \|x_{\text{src}}^{(n)} - g_{\text{target2src}}(g_{\text{src2target}}(x_{\text{src}}^{(n)}))\|_F$ and $\sum_n \|x_{\text{target}}^{(n)} - g_{\text{src2target}}(g_{\text{target2src}}(x_{\text{target}}^{(n)}))\|_F$ for $g_{\text{src2target}}$ and $g_{\text{target2src}}$



How Does Cycle GAN Work?

How Does Cycle GAN Work?

- Ideal: $g_{\text{src2target}}$ and $g_{\text{target2src}}$ learns to translate images

How Does Cycle GAN Work?

- Ideal: $g_{\text{src2target}}$ and $g_{\text{target2src}}$ learns to translate images
- Reality: $g_{\text{src2target}}$ and $g_{\text{target2src}}$ learns to hide information [6]



How Does Cycle GAN Work?

- Ideal: $g_{\text{src2target}}$ and $g_{\text{target2src}}$ learns to translate images
- Reality: $g_{\text{src2target}}$ and $g_{\text{target2src}}$ learns to hide information [6]



- *Unsupervised DNN models, including GANs, may not work as one may expect quality?*

Reference I

- [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou.
Wasserstein generative adversarial networks.
In *International Conference on Machine Learning*, pages 214–223, 2017.
- [2] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al.
Constitutional ai: Harmlessness from ai feedback.
arXiv preprint arXiv:2212.08073, 2022.
- [3] Andrew Brock, Theodore Lim, JM Ritchie, and Nick Weston.
Neural photo editing with introspective adversarial networks.
arXiv preprint arXiv:1609.07093, 2016.

Reference II

- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al.
Language models are few-shot learners.
Advances in neural information processing systems, 33:1877–1901, 2020.
- [5] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman.
Maskgit: Masked generative image transformer.
In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022.
- [6] Casey Chu, Andrey Zhmoginov, and Mark Sandler.
Cyclegan, a master of steganography.
arXiv preprint arXiv:1712.02950, 2017.

Reference III

- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding.
arXiv preprint arXiv:1810.04805, 2018.
- [8] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets.
In *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014.
- [9] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville.
Improved training of wasserstein gans.
arXiv preprint arXiv:1704.00028, 2017.

Reference IV

- [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel.
Denoising diffusion probabilistic models.
Advances in neural information processing systems, 33:6840–6851, 2020.
- [11] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.
Training compute-optimal large language models.
arXiv preprint arXiv:2203.15556, 2022.
- [12] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.
Lora: Low-rank adaptation of large language models.
arXiv preprint arXiv:2106.09685, 2021.

Reference V

- [13] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros.
Image-to-image translation with conditional adversarial networks.
arXiv preprint arXiv:1611.07004, 2016.
- [14] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen.
Progressive growing of gans for improved quality, stability, and variation.
arXiv preprint arXiv:1710.10196, 2017.
- [15] Diederik P Kingma and Max Welling.
Auto-encoding variational bayes.
arXiv preprint arXiv:1312.6114, 2013.
- [16] Durk P Kingma and Prafulla Dhariwal.
Glow: Generative flow with invertible 1x1 convolutions.
Advances in neural information processing systems, 31, 2018.

Reference VI

- [17] Quoc V Le and Tomas Mikolov.
Distributed representations of sentences and documents.
In *ICML*, volume 14, pages 1188–1196, 2014.
- [18] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al.
Photo-realistic single image super-resolution using a generative adversarial network.
arXiv preprint arXiv:1609.04802, 2016.
- [19] Daniel D Lee and H Sebastian Seung.
Learning the parts of objects by non-negative matrix factorization.
Nature, 401(6755):788–791, 1999.

Reference VII

- [20] Daniel D Lee and H Sebastian Seung.
Algorithms for non-negative matrix factorization.
In *Advances in neural information processing systems*, pages 556–562, 2001.
- [21] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al.
Retrieval-augmented generation for knowledge-intensive nlp tasks.
Advances in Neural Information Processing Systems, 33:9459–9474, 2020.
- [22] Calvin Luo.
Understanding diffusion models: A unified perspective.
arXiv preprint arXiv:2208.11970, 2022.

Reference VIII

- [23] Luke Metz, Ben Poole, David Pfau, and Jascha Sohl-Dickstein. Unrolled generative adversarial networks.
arXiv preprint arXiv:1611.02163, 2016.
- [24] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space.
arXiv preprint arXiv:1301.3781, 2013.
- [25] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality.
In *Advances in neural information processing systems*, pages 3111–3119, 2013.

Reference IX

- [26] Varun Nair, Elliot Schumacher, Geoffrey Tso, and Anitha Kannan. Dera: enhancing large language model completions with dialog-enabled resolving agents.
arXiv preprint arXiv:2303.17071, 2023.
- [27] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback.
arXiv preprint arXiv:2112.09332, 2021.
- [28] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback.
Advances in Neural Information Processing Systems, 35:27730–27744, 2022.

Reference X

- [29] Alec Radford, Luke Metz, and Soumith Chintala.
Unsupervised representation learning with deep convolutional
generative adversarial networks.
arXiv preprint arXiv:1511.06434, 2015.
- [30] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever,
et al.
Improving language understanding by generative pre-training.
OpenAI blog, 2018.
- [31] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei,
Ilya Sutskever, et al.
Language models are unsupervised multitask learners.
OpenAI blog, 2019.

Reference XI

- [32] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee.
Generative adversarial text to image synthesis.
arXiv preprint arXiv:1605.05396, 2016.
- [33] Salah Rifai, Pascal Vincent, Xavier Muller, Xavier Glorot, and Yoshua Bengio.
Contractive auto-encoders: Explicit invariance during feature extraction.
In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 833–840, 2011.
- [34] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton.
Dynamic routing between capsules.
In *Advances in Neural Information Processing Systems*, pages 3857–3867, 2017.

Reference XII

- [35] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen.
Improved techniques for training gans.
In *Advances in Neural Information Processing Systems*, pages 2226–2234, 2016.
- [36] Or Sharir, Barak Peleg, and Yoav Shoham.
The cost of training nlp models: A concise overview.
arXiv preprint arXiv:2004.08900, 2020.
- [37] Patrice Simard, Bernard Victorri, Yann LeCun, and John S Denker.
Tangent prop-a formalism for specifying selected invariances in an adaptive network.
In *NIPS*, volume 91, pages 895–903, 1991.

Reference XIII

- [38] Aaron van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al.
Conditional image generation with pixelcnn decoders.
In *Advances in Neural Information Processing Systems*, pages 4790–4798, 2016.
- [39] Aaron Van Oord, Nal Kalchbrenner, and Koray Kavukcuoglu.
Pixel recurrent neural networks.
In *International Conference on Machine Learning*, pages 1747–1756, 2016.
- [40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin.
Attention is all you need.
Advances in neural information processing systems, 30, 2017.

Reference XIV

- [41] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol.
Extracting and composing robust features with denoising autoencoders.
In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103. ACM, 2008.
- [42] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al.
Emergent abilities of large language models.
arXiv preprint arXiv:2206.07682, 2022.
- [43] Matthew D Zeiler and Rob Fergus.
Visualizing and understanding convolutional networks.
In *European conference on computer vision*, pages 818–833. Springer, 2014.

Reference XV

- [44] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. Generative visual manipulation on the natural image manifold. In *European Conference on Computer Vision*, pages 597–613. Springer, 2016.
- [45] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint arXiv:1703.10593*, 2017.