How to Al (Almost) Anything Lecture 10 – Modern Generative Al

Chanakya Ekbote

Multisensory Intelligence Group MIT Media Lab Adapted From: <u>https://diffusion.csail.mit.edu</u> Thanks to XKCD for the comics





Administrative Details

- 1. Reading Assignments due Tomorrow
 - 1. LLaMA-Adapter: Efficient Fine-tuning of Language Models with Zero-init Attention
 - 2. Cobra: Extending Mamba to Multi-Modal Large Language Model for Efficient Inference
- 2. Presentations are due this Thursday.
- 3. Hope you are all on track for your projects (talk to us in case of any blockers).

Todays Lecture

- 1. What are Generative Models?
- 2. Current State of the Art (Flow Matching)
- 3. Conditional Generation
- 4. Architectures
- 5. Tips to Train these Models

Generative Models





"Creating noise from data is easy; creating data from noise is generative modeling. Song et *al.* [30]



OpenAl GPT 40

OpenAl Sora



Prerequisites – Solving Differential Equations

$$\frac{dx}{dt} = t \qquad \longrightarrow \qquad \int dx = \int t dt \qquad \longrightarrow \qquad x = \frac{t^2}{2} + c$$



Prerequisites – Euler Integration

Trying to solve a differential equation like...





MODIFIED BAYES' THEOREM:

 $P(H|X) = P(H) \times \left(1 + P(C) \times \left(\frac{P(X|H)}{P(X)} - 1\right)\right)$

H: HYPOTHESIS X: OBSERVATION P(H): PRIOR PROBABILITY THAT H IS TRUE. P(X): PRIOR PROBABILITY OF OBSERVING X P(C): PROBABILITY THAT YOU'RE USING BAYESIAN STATISTICS CORRECTLY

PMF of a Single Coin Toss

Prerequisites - Probability

Outcome (0: Tails, 1: Heads)







 $x \sim N(\mu, \sigma^2)$ $\epsilon \sim N(0, 1)$ $x \sim \mu + \sigma \epsilon$

What is Generation?

Generation implies sampling from the data distribution



Why Sampling is a hard problem?

"Sampling the world exactly is impossible; approximations are everything."

Léon Bottou

1. We don't know *P*_{data}?

2. Even if we know, P_{data} we can't always sample from it.

Warmup

How topologists see the world:



Homeomorphic, bro.



Warmup



What if we know P_{data} ?



What if we know an approximation of P_{data} ?



Variational Autoencoders

μ Encoder Decoder \hat{x} X σ e d \odot $z = \mu + \sigma \odot \varepsilon$ sample N(0,I) -3 Minimize 2: $\frac{1}{2} \sum_{i=1}^{N} (\exp(\sigma_i) - (1 + \sigma_i) + {\mu_i}^2)$

Variational Autoencoders

Minimize 1: $(x - \hat{x})^2$

Denoising Diffusion Models



Diffusion Models: A Comprehensive Survey of Methods and Applications

What methods do we know?

Feature	VAE	Diffusion Models	Flow Matching			
Core Idea	Encode/decode with latent noise	Add noise and learn to reverse it	Learn a continuous flow from noise			
Training Objective	Minimize reconstruction + KL loss	Learn the score (gradient) of data	Match a vector field (OD based)			
Noise Handling	Noise in latent space	Progressive noise over time	Start from noise, smooth transform			
Sampling Speed	Very fast (one pass)	Slow (many denoising steps)	Faster (solving an ODE)			
Advantages	Simple, fast, interpretable	Very high-quality outputs	High quality + faster than diffusion			
Disadvantages	Blurry samples, limited expressiveness	Expensive, slow sampling	Newer, still developing			
Key Examples	VAE (2013) <i>,</i> β-VAE	DDPM, Stable Diffusion	Flow Matching (2023), Rectified Flow			

Flow Models

$$X_0 \sim p_{\text{init}}$$
$$\frac{\mathrm{d}}{\mathrm{d}t} X_t = u_t^{\theta}(X_t)$$

Algorithm 1 Sampling from a Flow Model with Euler method Require: Neural network vector field u_t^{θ} , number of steps n1: Set t = 02: Set step size $h = \frac{1}{n}$ 3: Draw a sample $X_0 \sim p_{\text{init}}$ 4: for i = 1, ..., n - 1 do 5: $X_{t+h} = X_t + hu_t^{\theta}(X_t)$ 6: Update $t \leftarrow t + h$ 7: end for 8: return X_1



Conditional and Marginal Probability Paths (Example)

$$\alpha_t = t, \beta_t = 1 - t$$

$$p_t(\cdot | z) = \mathcal{N}(\alpha_t z, \beta_t^2 I_d)$$

$$p_0(. | z) = \mathcal{N}(0, I_d)$$

$$p_1(. | z) = \delta_z$$



$$egin{aligned} & z \sim p_{ ext{data}}, \quad x \sim p_t(\cdot|z) \quad \Rightarrow x \sim p_t \ & p_t(x) = \int p_t(x|z) p_{ ext{data}}(z) \mathrm{d}z \ & p_0 = p_{ ext{init}} \quad ext{and} \quad p_1 = p_{ ext{data}}. \end{aligned}$$



 $z \sim p_{\text{data}}, \epsilon \sim p_{\text{init}} = \mathcal{N}(0, I_d) \Rightarrow x = \alpha_t z + \beta_t \epsilon \sim p_t$

Conditional Vector Field

$$X_0 \sim P_{init} \qquad \frac{dX_t}{dt} = u_t^{target}(X_t|z) \implies X_t \sim P_t(\cdot|z)$$

$$P_t(x_t|z) = \mathcal{N}(x_t; \alpha_t z, \beta_t^2 I_d) \implies u_t^{target}(x_t|z) = \left(\dot{\alpha}_t - \frac{\dot{\beta}_t}{\beta_t}\alpha_t\right)z + \frac{\dot{\beta}_t}{\beta_t}x_t$$



Conditional Vector Field (Intuition)

$$u_t^{target}(x_t|z) = \left(\dot{\alpha}_t - \frac{\dot{\beta}_t}{\beta_t}\alpha_t\right)z + \frac{\dot{\beta}_t}{\beta_t}x_t$$

			FI	ow N	/ato	:hing	g wi	ith Ti	rails	s (St	ep (D)					
			~ ~ ~	1	1	1	X	ź	ł	4	+	+	¥	¥	¥	1	1
	~~	~~	~ ~	~	\searrow	\searrow	\mathbf{x}	×	٩	٩	٠	+	¢	¢	*	¥	1
4 -		~ ~	~ ~	~	~	~	*	N	٩	٩	۴	۴	Þ	×	*	*	1
			~ ~	~	~	~	*	*	۹	٩	,	,	,	*	*	*	*
		$\rightarrow \rightarrow$		-	-	*	*	*	•	•	•	,	,	^	^	*	*
		$\rightarrow \rightarrow$	\rightarrow	-+	+	*	٠	*	*	•			^	4	4	*	*
2 -	$\rightarrow \rightarrow$	$\rightarrow \rightarrow$	$\rightarrow \rightarrow$	->	+	*	*	٣	٣		•			۲	٩	*	*
			·	-	+		-	٣	٣	*			`	٣	۳	-	-
			~ ~ ~			.#	*	*	*	4	4	b.	•	٣	*	*	Ψ.
~			~ ~ ~	~	1	π	1	#	*	4	4	k	•	*	*	*	*
0 -		~~~	~ ~	~	~	~	*	*	٩	٩	۲	7	,	,	1	1	1
			~ ~	~	*	*	*	*	٩	٩	,	,	,	*		#	×
				-	+	*	*	*	*	•		,	,	*		dir.	*
-2 -	$\rightarrow \rightarrow$	$\rightarrow \rightarrow$	\rightarrow	->	-	*	*	b-			•		-	4	-4	*	*
		$\rightarrow \rightarrow$	·	-	+	+	٠	٣	٣					٣	4	*	*
				-+	-#	~	~	*	*	*			`	٣	~	*	*
			~ ~ ~	-*	.*	ж	*	*	4	4	4	k	۲	۲	*	$\mathcal{T}_{\mathcal{T}}$	*
-4 -		~~	~ ~ ~	1	×	×	1	1	1	4	Å	k	k	x	×	*	*
			//	1	1	×	1	1	1	4	4	÷.	¥	۲	1	1	~
	//	//	//	1	1	1	1	1	1	1	4	ŧ	ł	Ł	X	X	1
	-	4	-	-2			Ċ)			ż				4		



Conditional Vector Field

$$P_t(x_t|z) = \mathcal{N}(x_t; \ \alpha_t z, \beta_t^2 I_d) \qquad \Longrightarrow \qquad u_t^{target}(x_t|z) = \left(\dot{\alpha}_t \ -\frac{\dot{\beta}_t}{\beta_t}\alpha_t\right) z \ + \frac{\dot{\beta}_t}{\beta_t} x_t$$

$$\alpha_t = t, \beta_t = 1 - t$$

$$P_t(x_t|z) = \mathcal{N}(x_t; tz, (1-t)^2 I_d) \implies u_t^{target}(x_t|z) = \frac{1+t}{1-t}z - \frac{1}{1-t}x_t$$

$$x_t = t z + (1 - t)\epsilon$$
 $u_t^{target}(x|z) = z - \epsilon$

Loss



$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{z \sim p_{\text{data}}, x \sim p_t(\cdot|z)} [\|u_t^{\theta}(x) - u_t^{\text{target}}(x|z)\|^2]$$



23

$$\mathcal{L}_{\mathrm{CFM}}(\theta) = \mathbb{E}_{z \sim p_{\mathrm{data}}, x \sim p_t(\cdot|z)} [\|u_t^{\theta}(x) - u_t^{\mathrm{target}}(x|z)\|^2].$$



Algorithm

Algorithm 3 Flow Matching Training Procedure (here for Gaussian CondOT path $p_t(x|z) = \mathcal{N}(tz, (1-t)^2)$)

Require: A dataset of samples $z \sim p_{data}$, neural network u_t^{θ}

- 1: for each mini-batch of data do
- 2: Sample a data example z from the dataset.
- 3: Sample a random time $t \sim \text{Unif}_{[0,1]}$.
- 4: Sample noise $\epsilon \sim \mathcal{N}(0, I_d)$
- 5: Set $x = tz + (1-t)\epsilon$
- 6: Compute loss

 $\mathcal{L}(\theta) = \|u_t^{\theta}(x) - (z - \epsilon)\|^2 \qquad (\text{General case:} = \|u_t^{\theta}(x) - u_t^{\text{target}}(x|z)\|^2)$

(General case: $x \sim p_t(\cdot|z)$)

7: Update the model parameters θ via gradient descent on $\mathcal{L}(\theta)$. 8: end for

Conditional Generation Loss

$$\mathcal{L}_{\mathrm{CFM}}^{\mathrm{guided}}(\theta) = \mathbb{E}_{(z,y) \sim p_{\mathrm{data}}(z,y), t \sim \mathrm{Unif}[0,1), x \sim p_t(\cdot|z)} \|u_t^{\theta}(x|y) - u_t^{\mathrm{target}}(x|z)\|^2$$

$$\mathcal{L}_{\rm CFM}^{\rm CFG}(\theta) = \mathbb{E}_{\Box} \| u_t^{\theta}(x|y) - u_t^{\rm target}(x|z) \|^2$$
$$\Box = (z, y) \sim p_{\rm data}(z, y), \ t \sim {\rm Unif}[0, 1), \ x \sim p_t(\cdot|z), \text{replace } y = \emptyset \text{ with prob. } \eta$$

$$\tilde{u}_t(x|y) = (1-w)u_t^{\text{target}}(x|\emptyset) + wu_t^{\text{target}}(x|y).$$





Metrics

Metric	What It Measures	Intuition	Good Value			
FID (Fréchet Inception Distance)	How close generated images are to real images (quality + diversity).	"How much the fake images <i>feel</i> like real ones."	Lower is better (e.g., FID < 10 is strong).			
CLIP Score	How well the generated image matches the text prompt.	"Did the model generate what I asked for?"	Higher is better.			
Precision / Recall	Precision = image realism. Recall = variety compared to real images.	"Are the images realistic (precision) and varied (recall)?"	High for both.			
Aesthetic Score	How beautiful or professional the images look.	"Would a human think this looks good?"	Higher is better.			

Model Architectures



Diffusion Transformer



Scalable Diffusion Models with Transformers

Conditioning with Language



Learning Transferable Visual Models From Natural Language Supervision

Diffusion in the Latent Space



Flow Matching in Latent Space

Stable Diffusion 3



Scaling Rectified Flow Transformers for High-Resolution Image Synthesis

Meta: MovieGen



Movie Gen: A Cast of Media Foundation Models

Tips to Train Flow Matching Models

- 1. Pick the right Loss Functions (L2 can be unstable)
- 2. Gradient Clipping maybe required
- 3. Use a Learning Rate Schedular (Liner Warmup + Cosine Decay)
- 4. Larger Batch Size is Preferred (Makes it more Stable)
- 5. Regularize

Tips to Debug Flow Matching Models

- 1. Visualize Trajectories if Possible (Test on Toy examples)
- 2. Does time (t) change anything?
- 3. Overfit on one sample and then on one batch
- 4. Check Loss Value on different examples
- 5. Check the norms of the output vector fields
- 6. Add small noise during validation (check if output changes wildly)

Further Reading

- 1. Flow Matching for Generative Modelling
- 2. Dreamcraft3d: Hierarchical 3d generation with bootstrapped diffusion prior
- 3. <u>Score-based generative modeling through stochastic differential equation</u>
- 4. <u>Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformer</u>
- 5. <u>Diffusion Models for Multi-Modal Generative Modelling</u>
- 6. <u>Multi-Track MusicLDM: Towards Versatile Music Generation with Latent Diffusion</u> <u>Model</u>
- 7. <u>https://lilianweng.github.io/posts/2021-07-11-diffusion-models/</u>
- 8. <u>https://yang-song.net/blog/</u>