

Deep Generative Models

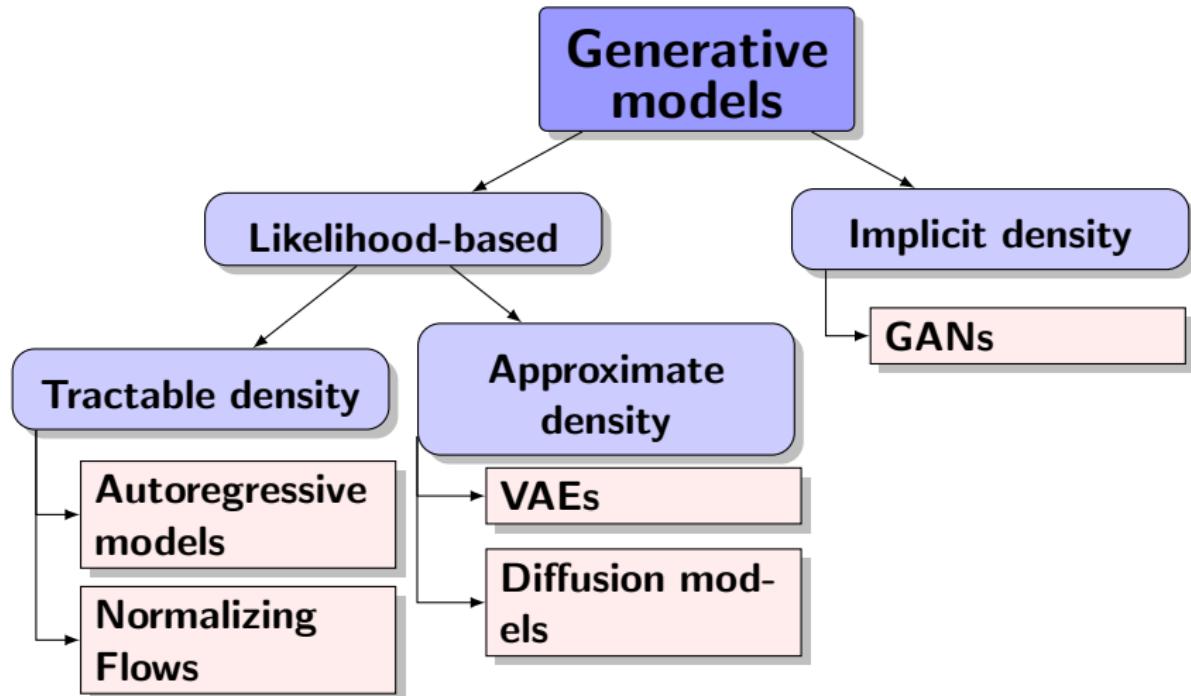
Lecture 1

Roman Isachenko



2025, Spring

A Zoo of Generative Models



Outline

1. Overview of Generative Models
2. Course Tricks
3. Problem Statement
4. Divergence Minimization Framework
5. Autoregressive Modeling

Outline

1. **Overview of Generative Models**
2. Course Tricks
3. Problem Statement
4. Divergence Minimization Framework
5. Autoregressive Modeling

VAE – The First Scalable Approach for Image Generation



DCGAN – The First Convolutional GAN for Image Generation



Radford A., Metz L., Chintala S. *Unsupervised representation learning with deep convolutional generative adversarial networks*, 2015

StyleGAN – High-Quality Facial Generation



Karras T., Laine S., Aila T. A style-based generator architecture for generative adversarial networks, 2018

Language Modeling at Scale

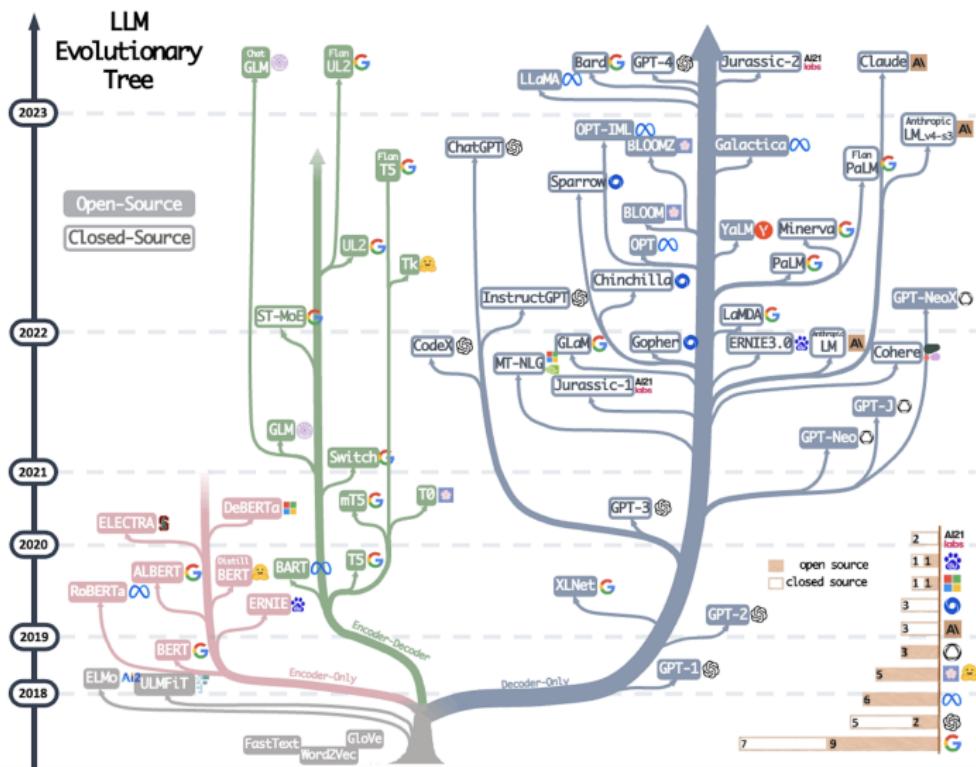
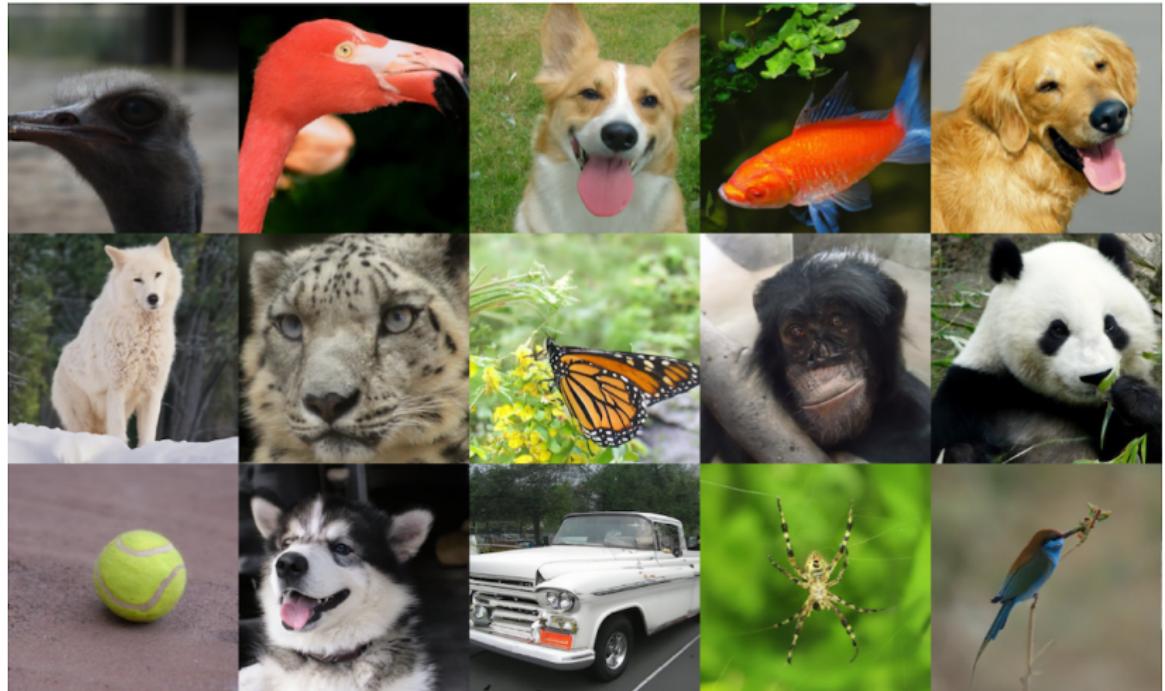


image credit:

<https://blog.biocomm.ai/2023/05/14/open-source-proliferation-lm-evolutionary-tree/>

Denoising Diffusion Probabilistic Models



Midjourney – Impressive Text-to-Image Results



image credit: <https://www.midjourney.com/explore>

Stable Diffusion 3 – Flow Matching



image credit: <https://stability.ai/news/stable-diffusion-3>

Sora – Video Generation



image credit: <https://openai.com/index/sora>

Outline

1. Overview of Generative Models
2. Course Tricks
3. Problem Statement
4. Divergence Minimization Framework
5. Autoregressive Modeling

Course Tricks (Part 1)

Log-Derivative Trick

Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be a differentiable function

$$\nabla \log f(\mathbf{x}) = \frac{1}{f(\mathbf{x})} \cdot \nabla f(\mathbf{x}).$$

Jensen's Inequality

Let $\mathbf{x} \in \mathbb{R}^m$ be a continuous random variable with a density $p(\mathbf{x})$ and $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be a convex function. Then

$$\mathbb{E}[f(\mathbf{x})] \geq f(\mathbb{E}[\mathbf{x}]).$$

Monte Carlo Estimation

Let $\mathbf{x} \in \mathbb{R}^m$ be a continuous random variable with a density $p(\mathbf{x})$ and let $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^d$ be a vector function. Then

$$\mathbb{E}_{p(\mathbf{x})}\mathbf{f}(\mathbf{x}) = \int p(\mathbf{x})\mathbf{f}(\mathbf{x})d\mathbf{x} \approx \frac{1}{n} \sum_{i=1}^n \mathbf{f}(\mathbf{x}_i), \quad \text{where } \mathbf{x}_i \sim p(\mathbf{x}).$$

Course Tricks (Part 2)

Change-of-Variable Theorem (CoV)

Let \mathbf{x} be a continuous random variable with a density $p(\mathbf{x})$ and $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^m$ be a differentiable, **invertible** function. If $\mathbf{y} = \mathbf{f}(\mathbf{x})$, then

$$p(\mathbf{y}) = p(\mathbf{x}) \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right) \right| = p(\mathbf{f}^{-1}(\mathbf{y})) \left| \det \left(\frac{\partial \mathbf{f}^{-1}(\mathbf{y})}{\partial \mathbf{y}} \right) \right|.$$

Proof (1D)

Assume f is a monotonically increasing function

$$F_Y(y) = P(Y \leq y) = P(x \leq f^{-1}(y)) = F_X(f^{-1}(y))$$

$$p(y) = \frac{dF_Y(y)}{dy} = \frac{dF_X(f^{-1}(y))}{dy} = \frac{dF_X(x)}{dx} \frac{df^{-1}(y)}{dy} = p(x) \frac{df^{-1}(y)}{dy}$$

Course Tricks (Part 3)

Law of the Unconscious Statistician (LOTUS)

Let $\mathbf{x} \in \mathbb{R}^m$ be a continuous random variable with a density $p(\mathbf{x})$ and let $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^m$ be a measurable function. If $\mathbf{y} = \mathbf{f}(\mathbf{x})$ then

$$\mathbb{E}_{p(\mathbf{y})}\mathbf{g}(\mathbf{y}) = \int p(\mathbf{y})\mathbf{g}(\mathbf{y})d\mathbf{y} = \int p(\mathbf{x})\mathbf{g}(\mathbf{f}(\mathbf{x}))d\mathbf{x} = \mathbb{E}_{p(\mathbf{x})}\mathbf{g}(\mathbf{f}(\mathbf{x})).$$

Dirac Delta Function

We can treat any deterministic variable \mathbf{x}_0 as a random variable with density $p(\mathbf{x}) = \delta(\mathbf{x} - \mathbf{x}_0)$.

$$\delta(\mathbf{x}) = \begin{cases} +\infty, & \mathbf{x} = 0; \\ 0, & \mathbf{x} \neq 0; \end{cases} \quad \int \delta(\mathbf{x})d\mathbf{x} = 1.$$

$$\mathbb{E}_{p(\mathbf{x})}\mathbf{f}(\mathbf{x}) = \int \delta(\mathbf{x} - \mathbf{x}_0)\mathbf{f}(\mathbf{x})d\mathbf{x} = \mathbf{f}(\mathbf{x}_0).$$

Outline

1. Overview of Generative Models
2. Course Tricks
3. Problem Statement
4. Divergence Minimization Framework
5. Autoregressive Modeling

Problem Statement

We are given i.i.d. samples $\{x_i\}_{i=1}^n \in \mathbb{R}^m$ from an unknown distribution $\pi(x)$.

Goal

We want to learn the distribution $\pi(x)$ to:

- ▶ evaluate $\pi(x)$ for new samples (i.e., how likely it is to obtain the object x) – density evaluation;
- ▶ generate new objects $x \sim \pi(x)$ – generation.

Challenge

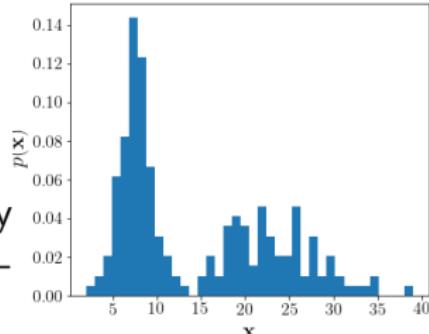
Data is complex and high-dimensional. For example, the dataset of images exists in $\mathbb{R}^{\text{width} \times \text{height} \times \text{channels}}$. The curse of dimensionality prevents us from knowing the exact density $\pi(x)$.

Histogram as a Generative Model

Let $x \sim \text{Categorical}(\pi)$. The histogram is fully determined by

$$\hat{\pi}_k = \hat{\pi}(x = k) = \frac{\sum_{i=1}^n [x_i = k]}{n}.$$

Problem: The curse of dimensionality (the number of bins grows exponentially).



MNIST example: 28x28 grayscale images, each image is $\mathbf{x} = (x_1, \dots, x_{784})$, where $x_i = \{0, 1\}$.

$$\pi(\mathbf{x}) = \pi(x_1) \cdot \pi(x_2|x_1) \cdot \dots \cdot \pi(x_m|x_{m-1}, \dots, x_1).$$

Hence, the histogram would have $2^{28 \times 28} - 1$ parameters to specify $\pi(\mathbf{x})$.

Question: How many parameters do we need in these cases?

$$\pi(\mathbf{x}) = \pi(x_1) \cdot \pi(x_2) \cdot \dots \cdot \pi(x_m);$$

$$\pi(\mathbf{x}) = \pi(x_1) \cdot \pi(x_2|x_1) \cdot \dots \cdot \pi(x_m|x_{m-1}).$$

Conditional Model

Practical Scenarios

In practice, it is common to build a conditional model $\pi(x|y)$.

- ▶ $y = \emptyset$, x – image \Rightarrow **unconditional image model**.
- ▶ y – class label, x – image \Rightarrow **class-conditional image model**.
- ▶ y – text prompt, x – image \Rightarrow **text-to-image model**.
- ▶ y – image, x – image \Rightarrow **image-to-image model**.
- ▶ y – image, x – text \Rightarrow **image captioning**.
- ▶ y – English text, x – Russian text \Rightarrow **sequence-to-sequence model (machine translation)**.
- ▶ y – sound, x – text \Rightarrow **speech-to-text model (automatic speech recognition)**.
- ▶ y – text, x – sound \Rightarrow **text-to-speech model**.

Outline

1. Overview of Generative Models
2. Course Tricks
3. Problem Statement
4. Divergence Minimization Framework
5. Autoregressive Modeling

Divergences

- ▶ Fix a probabilistic model $p(x|\theta)$ – a family of parameterized distributions.
- ▶ Instead of searching for the true $\pi(x)$ over all distributions, we learn a function approximation $p(x|\theta) \approx \pi(x)$.

What is a Divergence?

Let \mathcal{P} be the set of all possible probability distributions. Then $D : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$ is a **divergence** if

- ▶ $D(\pi||p) \geq 0$ for all $\pi, p \in \mathcal{P}$;
- ▶ $D(\pi||p) = 0$ if and only if $\pi \equiv p$.

Divergence Minimization Task

$$\min_{\theta} D(\pi||p),$$

where $\pi(x)$ is the true data distribution and $p(x|\theta)$ is the model distribution.

Forward KL vs. Reverse KL

Forward KL

$$KL(\pi||p) = \int \pi(\mathbf{x}) \log \frac{\pi(\mathbf{x})}{p(\mathbf{x}|\theta)} d\mathbf{x} \rightarrow \min_{\theta}$$

Reverse KL

$$KL(p||\pi) = \int p(\mathbf{x}|\theta) \log \frac{p(\mathbf{x}|\theta)}{\pi(\mathbf{x})} d\mathbf{x} \rightarrow \min_{\theta}$$

What is the difference between these two formulations?

Maximum Likelihood Estimation (MLE)

Let $\{\mathbf{x}_i\}_{i=1}^n$ be the set of given i.i.d. samples.

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^n p(\mathbf{x}_i|\theta) = \arg \max_{\theta} \sum_{i=1}^n \log p(\mathbf{x}_i|\theta).$$

Forward KL vs. Reverse KL

Forward KL

$$\begin{aligned} KL(\pi||p) &= \int \pi(\mathbf{x}) \log \frac{\pi(\mathbf{x})}{p(\mathbf{x}|\theta)} d\mathbf{x} \\ &= \int \pi(\mathbf{x}) \log \pi(\mathbf{x}) d\mathbf{x} - \int \pi(\mathbf{x}) \log p(\mathbf{x}|\theta) d\mathbf{x} \\ &= -\mathbb{E}_{\pi(\mathbf{x})} \log p(\mathbf{x}|\theta) + \text{const} \\ &\approx -\frac{1}{n} \sum_{i=1}^n \log p(\mathbf{x}_i|\theta) + \text{const} \rightarrow \min_{\theta}. \end{aligned}$$

Maximum likelihood estimation is equivalent to minimizing the Monte Carlo estimate of Forward KL.

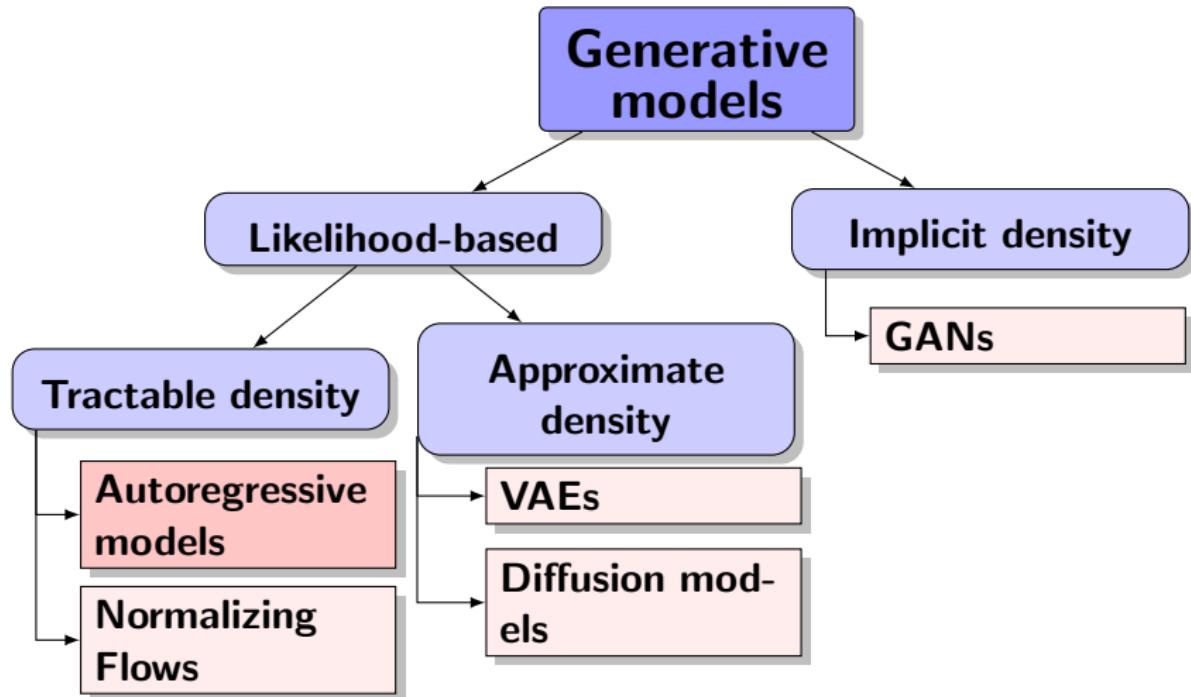
Reverse KL

$$\begin{aligned} KL(p||\pi) &= \int p(\mathbf{x}|\theta) \log \frac{p(\mathbf{x}|\theta)}{\pi(\mathbf{x})} d\mathbf{x} \\ &= \mathbb{E}_{p(\mathbf{x}|\theta)} [\log p(\mathbf{x}|\theta) - \log \pi(\mathbf{x})] \rightarrow \min_{\theta} \end{aligned}$$

Outline

1. Overview of Generative Models
2. Course Tricks
3. Problem Statement
4. Divergence Minimization Framework
5. Autoregressive Modeling

Generative Models Zoo (Revisited)



Autoregressive Modeling

MLE Problem

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^n p(\mathbf{x}_i | \theta) = \arg \max_{\theta} \sum_{i=1}^n \log p(\mathbf{x}_i | \theta).$$

- ▶ We aim to solve the problem using gradient-based optimization.
- ▶ We must be able to compute $\log p(\mathbf{x} | \theta)$ and its gradient $\frac{\partial \log p(\mathbf{x} | \theta)}{\partial \theta}$ efficiently.

Likelihood as Product of Conditionals

Let $\mathbf{x} = (x_1, \dots, x_m)$, and $\mathbf{x}_{1:j} = (x_1, \dots, x_j)$. Then

$$p(\mathbf{x} | \theta) = \prod_{j=1}^m p(x_j | \mathbf{x}_{1:j-1}, \theta); \quad \log p(\mathbf{x} | \theta) = \sum_{j=1}^m \log p(x_j | \mathbf{x}_{1:j-1}, \theta).$$

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^n \left[\sum_{j=1}^m \log p(x_{ij} | \mathbf{x}_{i,1:j-1}, \theta) \right].$$

Autoregressive Models

$$\log p(\mathbf{x}|\theta) = \sum_{j=1}^m \log p(x_j|\mathbf{x}_{1:j-1}, \theta).$$

- ▶ **Sampling is sequential:**
 - ▶ **sample** $\hat{x}_1 \sim p(x_1|\theta)$;
 - ▶ **sample** $\hat{x}_2 \sim p(x_2|\hat{x}_1, \theta)$;
 - ▶ \dots
 - ▶ **sample** $\hat{x}_m \sim p(x_m|\hat{\mathbf{x}}_{1:m-1}, \theta)$;
 - ▶ **the newly generated object is** $\hat{\mathbf{x}} = (\hat{x}_1, \hat{x}_2, \dots, \hat{x}_m)$.
- ▶ **Each conditional** $p(x_j|\mathbf{x}_{1:j-1}, \theta)$ **can be modeled by a neural network.**
- ▶ **Modeling all conditionals separately is infeasible. Shared parameters** θ **help avoid this.**

Autoregressive Models: MLP

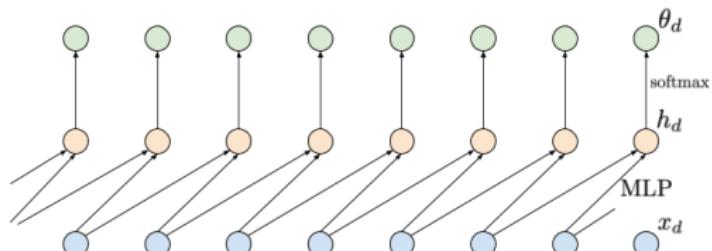
For large j , the conditional $p(x_j | \mathbf{x}_{1:j-1}, \theta)$ may be expensive.
Moreover, the history $\mathbf{x}_{1:j-1}$ has variable length.

Markov Assumption

$$p(x_j | \mathbf{x}_{1:j-1}, \theta) = p(x_j | \mathbf{x}_{j-d:j-1}, \theta), \quad d \text{ is a fixed model parameter.}$$

Example

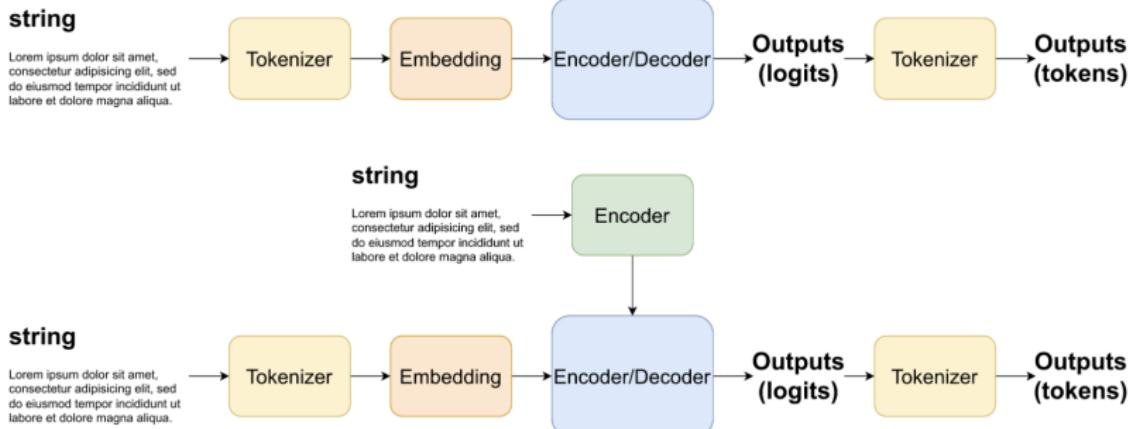
- ▶ $d = 2$;
- ▶ $x_j \in \{0, \dots, 255\}$;
- ▶ $\mathbf{h}_j = \text{MLP}_\theta(x_{j-1}, x_{j-2})$;
- ▶ $\pi_j = \text{softmax}(\mathbf{h}_j)$;
- ▶ $p(x_j | x_{j-1}, x_{j-2}, \theta) = \text{Categorical}(\pi_j)$.



We can also model continuous distributions instead of discrete ones.

Autoregressive Models: LLM

$$p(x_j | \mathbf{x}_{1:j-1}, \theta) = p(x_j | \mathbf{x}_{j-d:j-1}, \theta), \quad d \text{ is a fixed context length.}$$

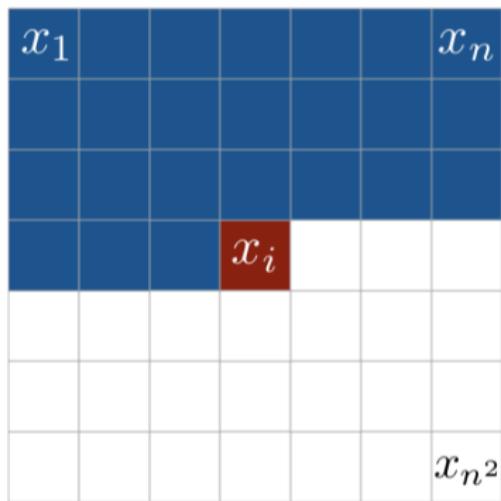


Autoregressive Models for Images

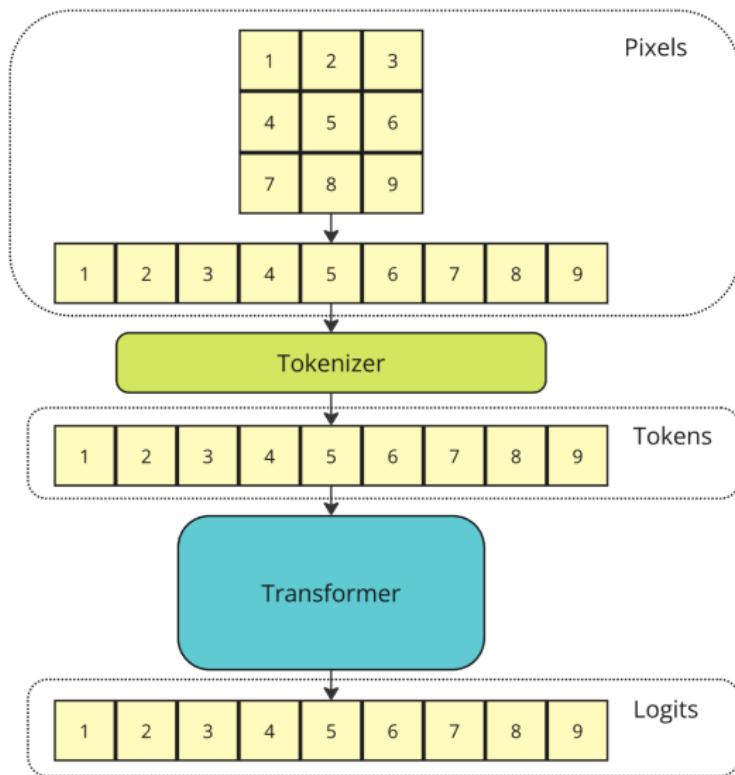
How do we model the distribution $\pi(\mathbf{x})$ of natural images?

$$p(\mathbf{x}|\theta) = \prod_{j=1}^{\text{width} \times \text{height}} p(x_j | \mathbf{x}_{1:j-1}, \theta).$$

- ▶ We must specify an order of the image pixels.
Raster ordering is a common approach.
- ▶ We can also incorporate RGB channel dependencies.



Autoregressive Models: ImageGPT



Summary

- ▶ We aim to approximate the data distribution for both density estimation and sample generation.
- ▶ Fitting a model distribution to the true data distribution can be posed as divergence minimization.
- ▶ Minimizing the forward KL is equivalent to maximizing the likelihood (MLE).
- ▶ Autoregressive models decompose the distribution into sequential conditionals.
- ▶ Sampling from autoregressive models is straightforward, but it must proceed sequentially.
- ▶ To compute the density, we multiply all conditionals $p(x_j | \mathbf{x}_{1:j-1}, \theta)$.
- ▶ ImageGPT flattens the image (raster ordering) so that the Transformer model can process it.