

# Deep Generative Models

## Lecture 1

Roman Isachenko



Autumn, 2022

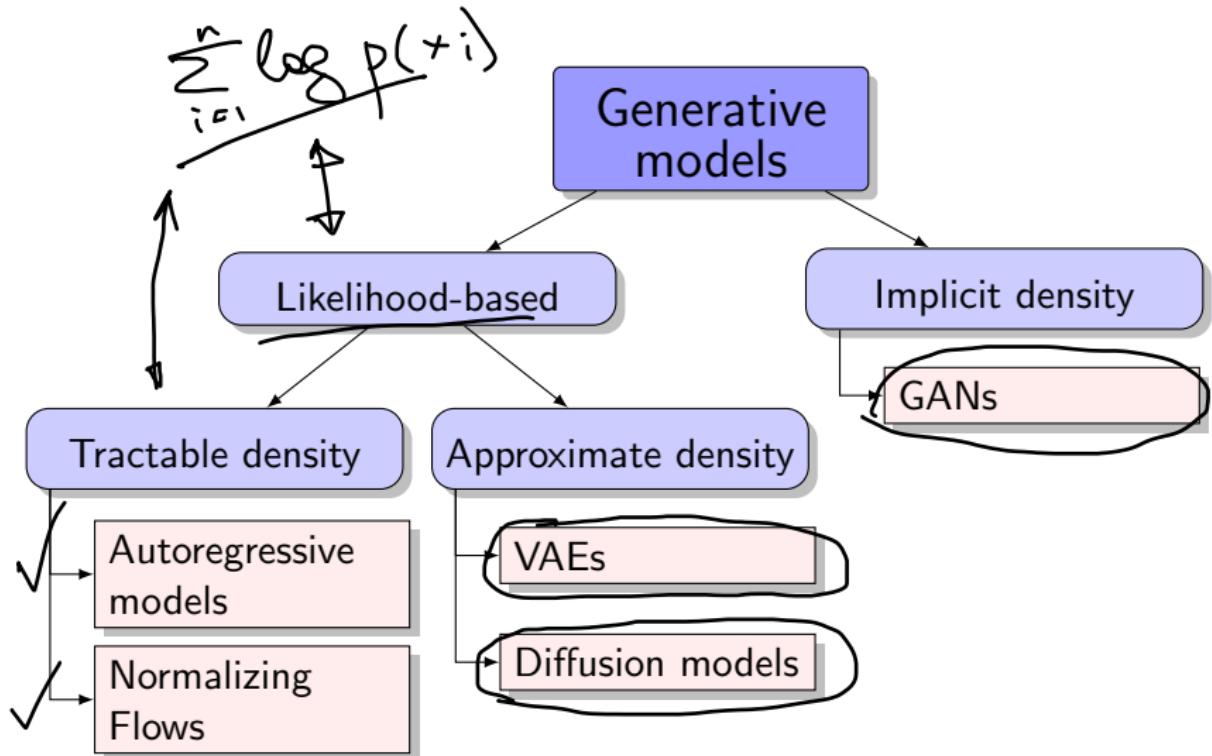
# Outline

1. Generative models overview
2. Problem statement
3. Divergence minimization framework
4. Autoregressive modelling

# Outline

1. Generative models overview
2. Problem statement
3. Divergence minimization framework
4. Autoregressive modelling

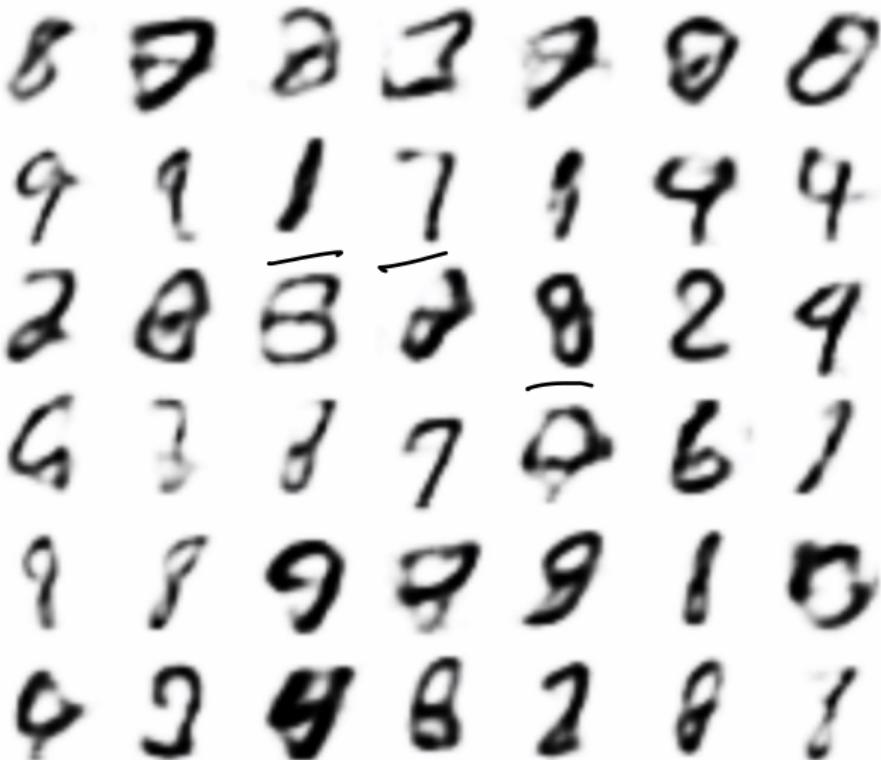
# Generative models zoo



## VAE – first scalable approach for image generation

MNIST

Gok



# DCGAN – first convolutional GAN for image generation

64



64



LSunBedrooms



# StyleGAN – high quality generation of faces



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

## VQ-VAE-2 – high quality generation without GANs

---

1024

1024



# Language modelling at scale

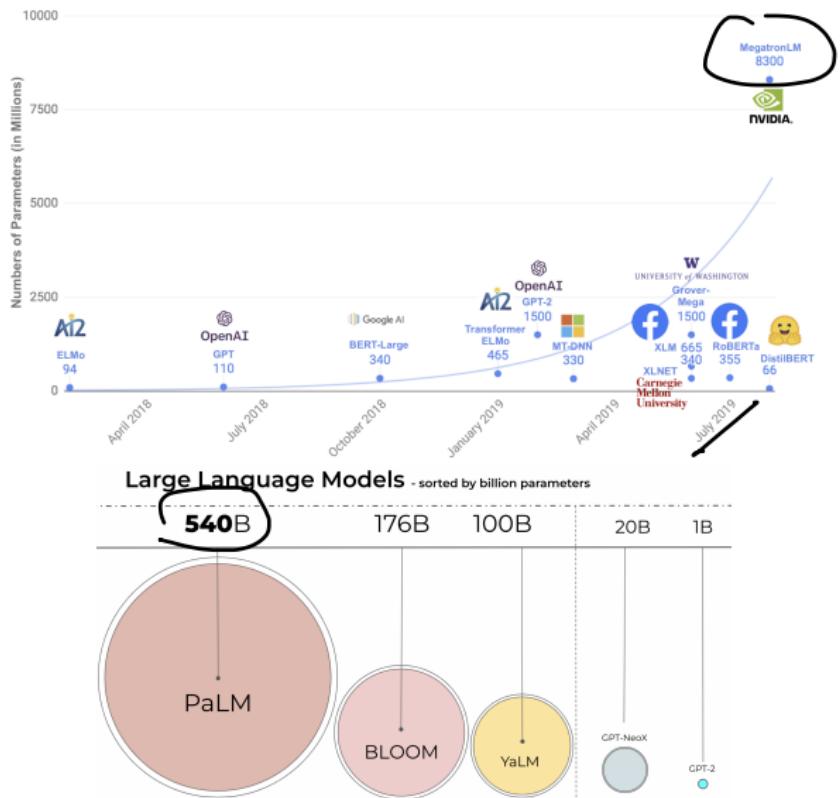


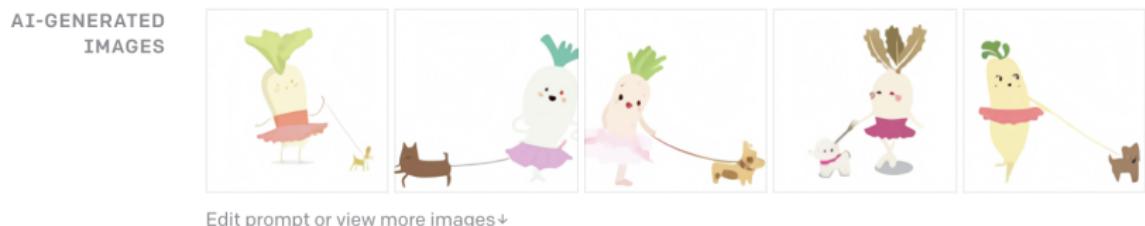
image credit: <http://jalammar.github.io/illustrated-gpt2>

image credit: <https://huggingface.co/blog/hf-bitsandbytes-integration>

# DALL-E – cross-modal image-text model



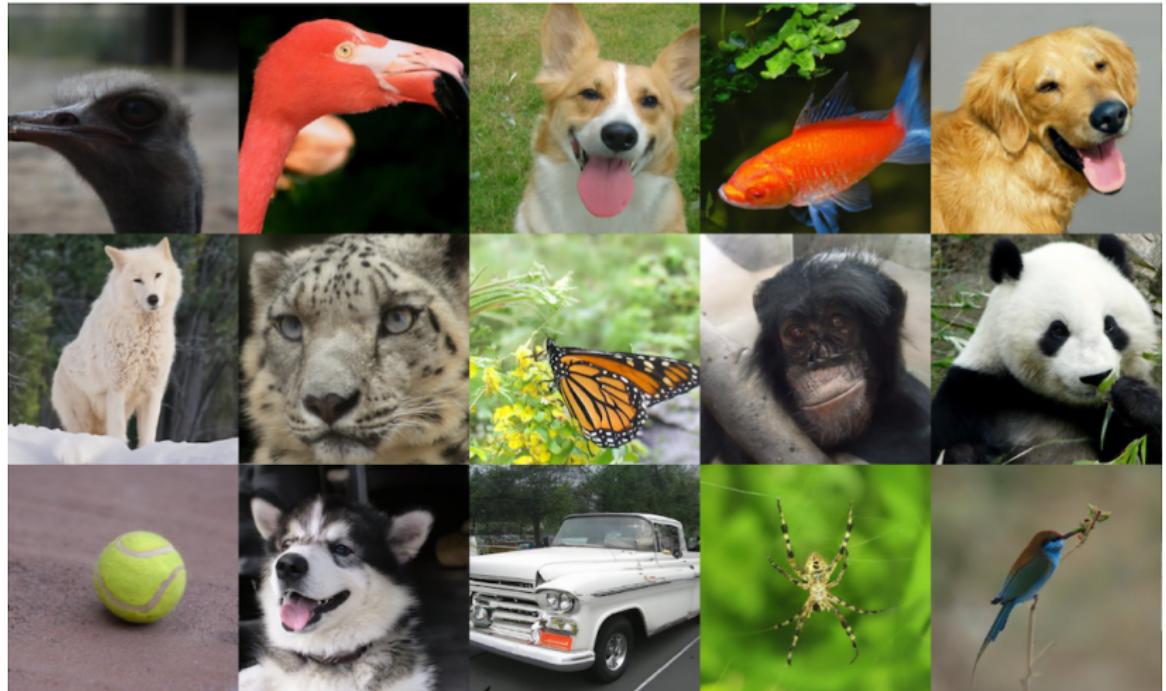
TEXT PROMPT    an illustration of a baby daikon radish in a tutu walking a dog



---

image credit: <https://openai.com/blog/dall-e/>  
Ramesh A. et al. Zero-shot text-to-image generation, 2021

# DDPM - diffusion model



# Stable Diffusion - awesome text to image results



---

Rombach R., et al. *High-Resolution Image Synthesis with Latent Diffusion Models*,  
2021

<https://github.com/CompVis/stable-diffusion>  
LAION-5B dataset

# Outline

1. Generative models overview
2. Problem statement
3. Divergence minimization framework
4. Autoregressive modelling

## Course tricks 1

$$f: \mathbb{R}^m \rightarrow \mathbb{R}^k$$
$$x \in \mathbb{R}^m$$

### Law of the unconscious statistician (LOTUS)

Let  $\underline{x}$  be a random variable with density  $p_x(\underline{x})$  and let  $\underline{y} = f(\underline{x})$  with density  $p_y(\underline{y})$ . Then

$$\mathbb{E}_{p_y} \underline{y} = \int p_y(\underline{y}) \underline{y} d\underline{y} = \int p_x(\underline{x}) f(\underline{x}) d\underline{x} = \mathbb{E}_{p_x} f(\underline{x}).$$

### Monte-Carlo estimation

Expected value could be estimated using only the samples:

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

$$x \sim P(+)$$
$$P(+) = 0.2$$

## Course tricks 2



### Jensen's Inequality

Let  $\underline{x}$  be a random variable and  $f(\cdot)$  is a convex function. Then

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

$$P(A \cap B) \geq P(A|B)P(B)$$

### Decomposition to conditionals

Let  $\underline{x} = (x_1, x_2, \dots, x_m) \in \mathbb{R}^m$  be a random variable. Then

$$p(\underline{x}) = \underbrace{p(x_1)}_{\text{ }} \cdot \underbrace{p(x_2|x_1)}_{\text{ }} \cdot \underbrace{p(x_3|x_2, x_1)}_{\text{ }} \cdot \dots \cdot \underbrace{p(x_m|x_{m-1}, \dots, x_1)}_{\text{ }}$$

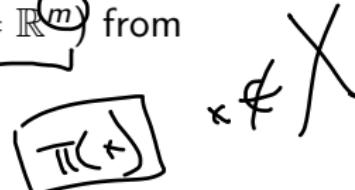
$$p(x_2|x_1) = p(x_2)$$

$$p(\underline{x}) = \prod_{i=1}^m p(x_i)$$

## Problem statement



We are given i.i.d. samples  $\{\mathbf{x}_i\}_{i=1}^n \in \mathcal{X}$  (e.g.  $\mathcal{X} = \mathbb{R}^m$ ) from unknown distribution  $\pi(\mathbf{x})$ .



## Goal

We would like to learn a distribution  $\pi(\mathbf{x})$  for

- ▶ evaluating  $\pi(\mathbf{x})$  for new samples (how likely to get object  $\mathbf{x}$ ?);
- ▶ sampling from  $\pi(\mathbf{x})$  (to get new objects  $\mathbf{x} \sim \pi(\mathbf{x})$ ).

## Challenge

Data is complex and high-dimensional. E.g. the dataset of images lies in the space  $\mathcal{X} \subset \mathbb{R}^{\text{width} \times \text{height} \times \text{channels}}$ .

$$\mathbf{x} \in \mathbb{R}^m \quad m \sim 10^8$$

## Histogram as a generative model

⊕

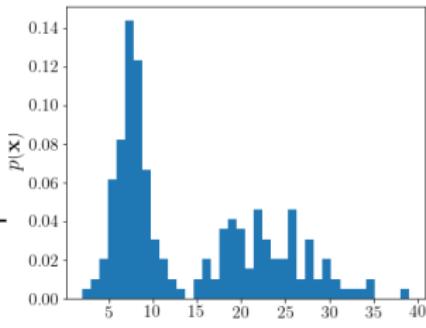
$$x \in \{1, \dots, K\}$$

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

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

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

$$x \in \text{Be}(p) \subset \{0, 1\}$$



**MNIST example:** 28x28 gray-scaled images, each image is  $x = (x_1, \dots, x_{784})$ , where  $x_i \in \{0, 1\} \subset \text{Be}(p)$

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

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

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

$$\frac{\pi(x) = \pi(x_1) \cdot \pi(x_2) \cdot \dots \cdot \pi(x_m); m}{\pi(x) = \pi(x_1) \cdot \pi(x_2|x_1) \cdot \dots \cdot \pi(x_m|x_{m-1})} = 2^{m-1}$$

# Outline

1. Generative models overview
2. Problem statement
3. Divergence minimization framework
4. Autoregressive modelling

# Divergences

Fix probabilistic model  $p(\mathbf{x}|\theta)$  – the set of parameterized distributions.

Instead of searching true  $\pi(\mathbf{x})$  over all probability distributions, learn function approximation  $p(\mathbf{x}|\theta) \approx \pi(\mathbf{x})$ .

## What is a divergence?

Let  $\mathcal{S}$  be the set of all possible probability distributions. Then

$D : \overline{\mathcal{S}} \times \mathcal{S} \rightarrow \mathbb{R}_+$  is a divergence if

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

• comm exp.  
• need to try?.

## Divergence minimization task

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

where  $\pi(\mathbf{x})$  is a true data distribution,  $p(\mathbf{x}|\theta)$  is a model distribution.

# f-divergence family

## f-divergence

$$D_f(\pi || p) = \mathbb{E}_{p(x)} f\left(\frac{\pi(x)}{p(x)}\right) = \int p(x) f\left(\frac{\pi(x)}{p(x)}\right) dx.$$

Here  $f : \mathbb{R}_+ \rightarrow \mathbb{R}$  is a convex, lower semicontinuous function satisfying  $f(1) = 0$ .

Name	$D_f(P  Q)$	Generator $f(u)$
Kullback-Leibler	$\int p(x) \log \frac{p(x)}{q(x)} dx$	$u \log u$
Reverse KL	$\int q(x) \log \frac{q(x)}{p(x)} dx$	$- \log u$
Pearson $\chi^2$	$\int \frac{(q(x)-p(x))^2}{p(x)} dx$	$(u-1)^2$
Squared Hellinger	$\int \left( \sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx$	$(\sqrt{u}-1)^2$
Jensen-Shannon	$\frac{1}{2} \int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx$	$-(u+1) \log \frac{1+u}{2} + u \log u$
GAN	$\int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx - \log(4)$	$u \log u - (u+1) \log(u+1)$

# Forward KL vs Reverse KL

## Forward KL

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

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

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

## Forward KL vs Reverse KL

### Forward KL

$$\log \frac{a}{b} = \log a - \log b$$

$$\begin{aligned} \min_{\theta} & \leftarrow KL(\pi || p) = \int \pi(\mathbf{x}) \log \frac{\pi(\mathbf{x})}{p(\mathbf{x}|\theta)} d\mathbf{x} \\ &= \int \pi(\mathbf{x}) \log \pi(\mathbf{x}) d\mathbf{x} - \int \pi(\mathbf{x}) \log p(\mathbf{x}|\theta) d\mathbf{x} \\ &= -\mathbb{E}_{\pi(\mathbf{x})} \log p(\mathbf{x}|\theta) + \text{const}(\theta) \end{aligned}$$

$x_i \sim \pi(\mathbf{x})$

**MC**  $\approx \left[ \frac{1}{n} \sum_{i=1}^n \log p(\mathbf{x}_i|\theta) + \text{const} \right] \xrightarrow{\max_{\theta}} \min_{\theta}$

Maximum likelihood estimation is equivalent to minimization of the Monte-Carlo estimate of forward KL.

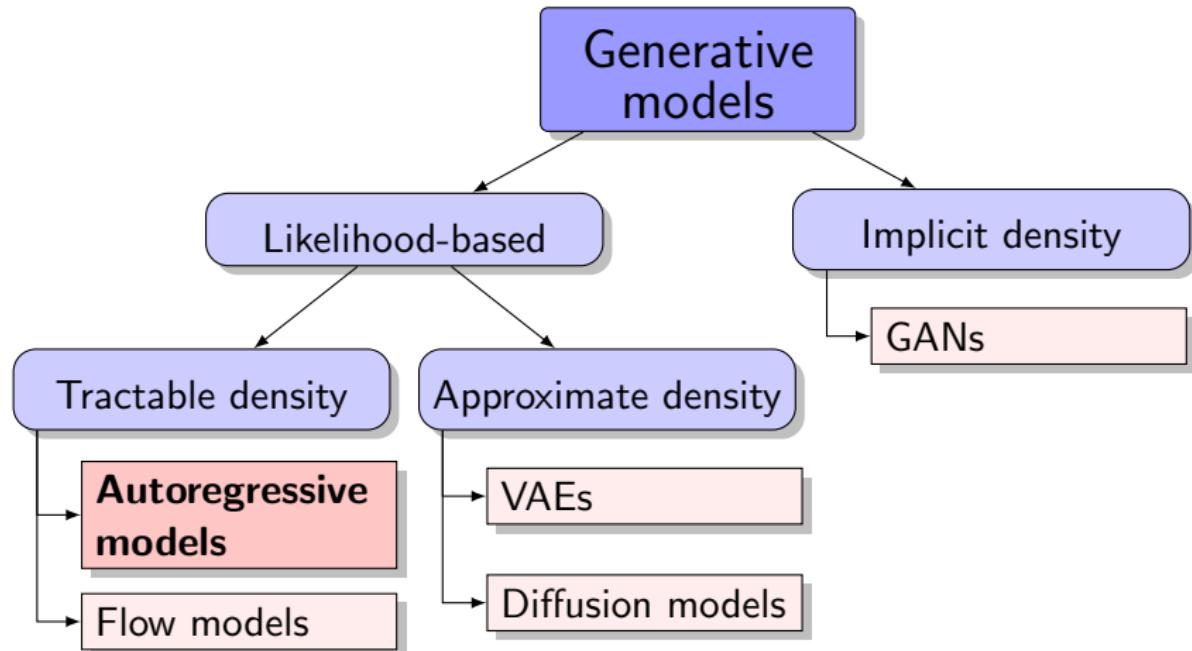
### Reverse KL

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

# Outline

1. Generative models overview
2. Problem statement
3. Divergence minimization framework
4. Autoregressive modelling

# Generative models zoo



Autoregressive modelling  $\oplus$   $\sum_{i=1}^n \sum_{j=1}^m \log p(x_{i:j} | x_{1:j-1}^i)$

MLE problem  $\equiv$   $\min_{\theta} -n \text{FKL}(\theta)$

$$\hat{\theta}^* = \arg \max_{\theta} p(\mathbf{X}|\theta) = \arg \max_{\theta} \prod_{i=1}^n p(x_i|\theta) = \arg \max_{\theta} \underbrace{\sum_{i=1}^n \log p(x_i|\theta)}_{\text{FKL}}$$

- We would like to solve the problem using gradient-based optimization.
- We have to efficiently compute  $\log p(\mathbf{x}|\theta)$  and  $\frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta}$ .

Likelihood as product of conditionals  $x_{1:n} = \{x_1\}$

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

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

**Example:**  $p(x_1, x_2, x_3) = p(x_1) \cdot p(x_2 | x_1) \cdot p(x_3 | x_1, x_2)$ .

## Autoregressive models

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

- ▶ Sampling is sequential:
  - ▶ sample  $\hat{x}_1 \sim p(x_1 | \theta)$ ;
  - ▶ sample  $\hat{x}_2 \sim p(x_2 | \hat{x}_1, \theta)$ ;
  - ▶ ...
  - ▶ sample  $\hat{x}_m \sim p(x_m | \hat{x}_{1:m-1}, \theta)$ ;
  - ▶ new generated object is  $\hat{\mathbf{x}} = (\hat{x}_1, \hat{x}_2, \dots, \hat{x}_m)$ .
- ▶ Each conditional  $p(x_j | \underline{\mathbf{x}_{1:j-1}}, \theta)$  could be modelled by neural network.
- ▶ Modelling all conditional distributions separately is infeasible and we would obtain separate models. To extend to high dimensions we could share parameters  $\theta$  across conditionals.

# Autoregressive models



For large  $j$  the conditional distribution  $p(x_j | \underline{x}_{1:j-1}, \theta)$  could be infeasible. Moreover, the history  $\underline{x}_{1:j-1}$  has non-fixed length.

## Markov assumption

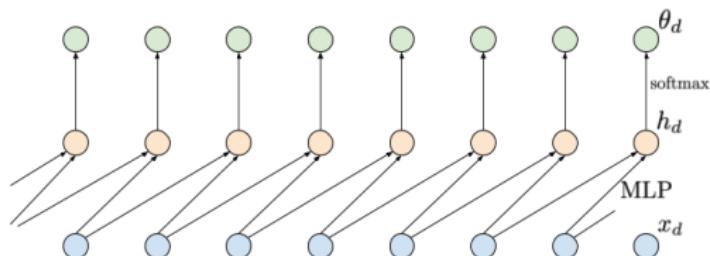
$$p(x_j | \underline{x}_{1:j-1}, \theta) = p(x_j | \underline{x}_{j-d:j-1}, \theta), \quad d \text{ is a fixed model parameter.}$$

## Example

$$\begin{array}{l} \blacktriangleright d = 2; \\ \blacktriangleright x_j \in \mathbb{R} \end{array}$$



- ▶  $d = 2$ ;
- ▶ ~~history of length  $d$~~
- ▶  $\mathbf{h}_j = \text{MLP}_\theta(x_{j-1}, x_{j-2});$
- ▶  $\pi_j = \text{softmax}(\mathbf{h}_j);$

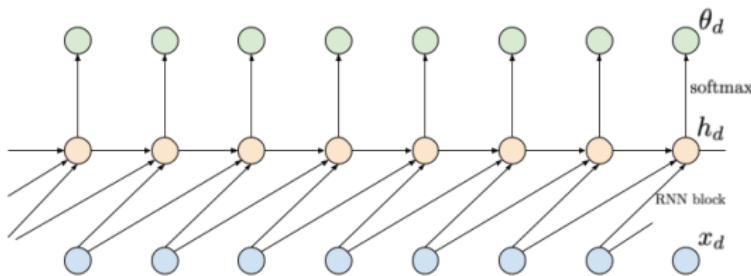


- ▶  $p(x_j | x_{j-1}, x_{j-2}, \theta) = \text{Categorical}(\pi_j).$
- Is it possible to model continuous distributions instead of discrete one?
- $N(\mu, \sigma^2)$        $\mu(x_{j-1}, x_{j-2})$        $\sigma^2(\dots)$

## Autoregressive models

- ▶ Previous model has **limited** memory  $d$ . It is insufficient for many modalities (e.g. for images and text).
- ▶ Recurrent NN fixes this problem and potentially could learn long-range dependencies:

$$p(x_j | \mathbf{x}_{1:j-1}, \theta) = p(x_j | \mathbf{h}_j, \theta), \quad \mathbf{h}_j = \text{RNN}(\mathbf{x}_{j-d:j-1}, \mathbf{h}_{j-1})$$



- ▶ Sequential computation of all conditionals  $p(x_j | \mathbf{x}_{1:j-1}, \theta)$ , hence, the training is slow.
- ▶ RNN suffers from vanishing and exploding gradients.

## Summary

- ▶ We are trying to approximate the distribution of samples for density estimation and generation of new samples.
- ▶ To fit model distribution to the real data distribution one could use divergence minimization framework.
- ▶ Minimization of forward KL is equivalent to the MLE problem.
- ▶ Autoregressive models decompose the distribution to the sequence of the conditionals.
- ▶ Sampling from the autoregressive models is trivial, but sequential
  - ▶ sample  $\hat{x}_1 \sim p(x_1)$ ;
  - ▶ sample  $\hat{x}_2 \sim p(x_2 | \hat{x}_1)$ ;
  - ▶ ....
- ▶ Density estimation:

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

- ▶ Autoregressive models work on both continuous and discrete data.