# Generative IA

Score matching, Langevin dynamics, diffusion models

Arnaud Breloy   `arnaud.breloy@lecnam.net`

15 décembre 2024

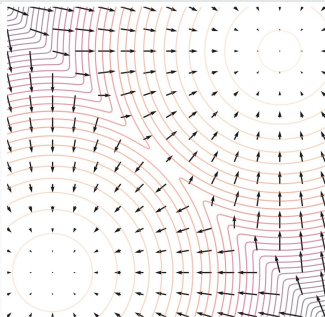Conservatoire national des arts & métiers

## Score Function

### Definition

The score function for a given distribution $p$ is given by

$$s(x) = \nabla_x \log p(x)$$

- Give the direction of the distributions models
- Can be used for sampling

## Score based models

### Definition

Given a dataset $\{x_1, ..., x_N\} \in \mathcal{X}$ with $\forall i, x_i \sim p$, a score based model $s_\theta$ is learned to retrive the score of the distribution :

$$s_\theta(x) \simeq \nabla_x \log p(x)$$

$$\theta = \arg\min_\theta \mathbb{E}_{x \sim p}\left[||\nabla_x \log p(x) - s_\theta(x)||_2^2\right]$$

$\hookrightarrow$ Requires # dimension($\mathcal{X}$) back-propagation
$\hookrightarrow$ Untractable in high-dimension

### Link with energy-based models

Key feature : normalization constant not required

$$\begin{aligned}
\nabla_x \log p_\theta(x) &= \nabla_x \log \frac{e^{-f_\theta(x)}}{Z_\theta} \\
&= \nabla_x \log e^{-f_\theta(x)} + \nabla_x \log(Z_\theta) \\
&= -\nabla_x f_\theta(x)
\end{aligned}$$

# Score based models

## Definition

Given a dataset $\{x_1, ..., x_N\} \in \mathcal{X}$ with $\forall i$, $x_i \sim p$, a score based model $s_\theta$ is learned to retrive the score of the distribution :

$$s_\theta(x) \simeq \nabla_x \log p(x)$$

$$\theta = \arg\min_\theta \mathbb{E}_{x \sim p}\left[||\nabla_x \log p(x) - s_\theta(x)||_2^2\right]$$

$\hookrightarrow$ Requires # dimension$(\mathcal{X})$ back-propagation
$\hookrightarrow$ Untractable in high-dimension

## Link with energy-based models

Key feature : normalization constant not required

$$\begin{aligned}
\nabla_x \log p_\theta(x) &= \nabla_x \log \frac{e^{-f_\theta(x)}}{Z_\theta} \\
&= \nabla_x \log e^{-f_\theta(x)} + \nabla_x \log(Z_\theta) \\
&= -\nabla_x f_\theta(x)
\end{aligned}$$

# Score approximation techniques

## Sliced-score matching

Project the score on random directions $v \sim p_v$ before minimazing the loss

$$\mathbb{E}_{v \sim p_v} \mathbb{E}_{x \sim p} \left( v^T \nabla_x \log p(x) - v^T s_\theta(x) \right)^2$$

## Denoising score matching

Equivalence between denoising autoencoder (DAE) objective and score matching

- Matching the score of a noise perturbed distribution

$$\mathbb{E}_{\tilde{x} \sim q_\sigma(\tilde{x})} || \nabla_{\tilde{x}} \log q_\sigma(\tilde{x}) - s_\theta(\tilde{x}) ||_2^2$$

- Equivalent to denoising score matching $\Leftrightarrow$ DAE objective

$$\mathbb{E}_{\sim p} \mathbb{E}_{\tilde{x} \sim q_\sigma(\tilde{x}|x)} || \underbrace{\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x)}_{\text{Tractable}} - s_\theta(\tilde{x}) ||_2^2$$

# Denoising Score Matching

## Training

- Sample a batch of data $\{x_1, ..., x_n\} \sim p(x)$
- Sample noisy data $\{x_1, ..., x_n\} \sim q_\sigma(\tilde{x}|x)$
- Estimate the denoising score loss :

$$\frac{1}{n} \sum_{i=1}^{n} ||s_\theta(\tilde{x}_i) - \nabla_{\tilde{x}} \log q_\sigma(\tilde{x}_i|x_i)||_2^2$$

- In the case of additive Gaussian noise :

$$\frac{1}{n} \sum_{i=1}^{n} ||s_\theta(\tilde{x}_i) - \frac{x_i - \tilde{x}_i}{\sigma^2}||_2^2$$

- Compute gradient descent
- $\sigma$ must be small

## Langevin Dynamic Sampling

### Stochastic sampling process

- $x_0 \sim \pi(x) \leftarrow$ random initialization
- repeat for $t \in 1, 2, ..., T$

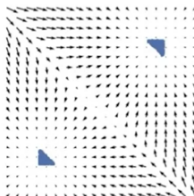$$x_t \leftarrow x_t + \frac{\epsilon}{2} \nabla_x \log p(x_{t-1}) + \epsilon z_t$$

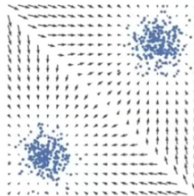  with $\epsilon << 0$ and $T \rightarrow \infty$.

- Using the score estimation :

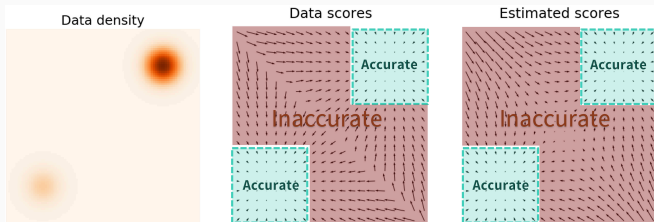$$x_t \leftarrow x_t + \frac{\epsilon}{2} s_\theta(x) + \epsilon z_t$$



Score function          Follow the scores          Follow the noisy scores

**Problem** Estimated scores are only accurate in high density regions.
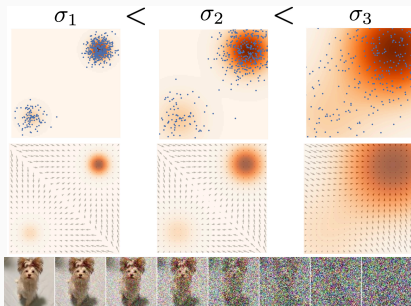
## Score matching with multiple noise scale

Adding noise to the data spread their distribution filling low density regions. The score is then learned on noisy data perturbed with different noise scale :

$$\mathcal{L}(\theta) = \sum_{l=1}^{L} \lambda(i) \mathbb{E}_{x \sim p_{\sigma_i}} ||\nabla_x \log p_{\sigma_i}(x) - s_\theta(x, \sigma_i)||_2^2,$$
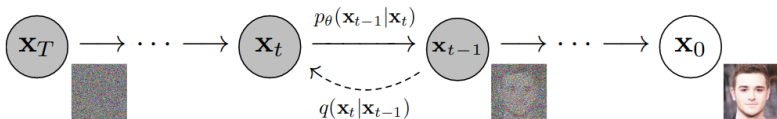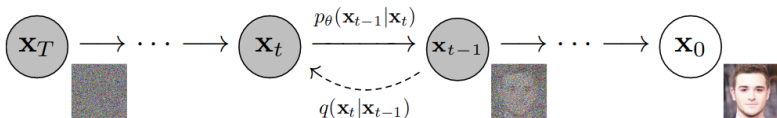
with $L$ the number of noise scales.

# Diffusion Models

Different formulation equivalent to score based approaches

- Noise is progressively added to the image
- The model learn a denosier able to retrieve $x_{t-1}$ from $x_t$
- Sampling is then performed by following $T$ denoising step starting from $x_T \sim \mathcal{N}(0, \sigma^2 I)$

# Forward and Backward Processes



## Forward Process

Gradually adds noise to the image. It is deifned as a Markow chain :

- $q(x_1, ..., x_T | x_0) = \prod_{t=1}^{T} q_\theta(x_t | x_{t-1})$
- $q(x_t | x_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t x_{t-1}}, \beta_t I)$

## Backward Process

Gradually substract the noise. It is deifned as a Markow chain :

- $p(x_T) = \mathcal{N}(0, \sigma^1 I)$
- $p_\theta(x_0, ..., x_T) = p(x_T) \prod_{t=1}^{T} p_\theta(x_{t-1} | x_t)$
- $p_\theta(x_{t-1} | x_t) = \mathcal{N}(\boldsymbol{\mu}_\theta, \boldsymbol{\Sigma}_\theta(x_t, t))$

## Evidencial Lower Bound

### Loss

Training is performed by minmizing the following bound on the negative likelihood :

$$\mathbb{E}[-\log p_\theta(x_0)] \leq \mathbb{E}_q\left[-\log\frac{p_\theta(x_0,...,x_T)}{q(x_1,...,x_T|x_0)}\right]$$

$$= \mathbb{E}_q\left[\text{constant} + \sum_{t>1}\text{KL}(q(x_{t-1}|x_t,x_0)||p_\theta(x_{t-1}|x_t))\right.$$

$$\left.- \log p_\theta(x_0|x_1)\right],$$

with $q(x_{t-1}|x_t,x_0) = \mathcal{N}(\tilde{\boldsymbol{\mu}}_t(x_t,x_0),\tilde{\beta}_t I)$

### Training

Training is performed by sampling a batch of time-steps and minimizing the sum with respect to it. The Kullback-Leiber divergences can be written as

$$\text{KL}(q(x_{t-1}|x_t,x_0)||p_\theta(x_{t-1}|x_t)) = \frac{1}{2\sigma_t^2}||\tilde{\boldsymbol{\mu}}_t(x_t,x_0) - \boldsymbol{\mu}_\theta(x_t,t)||_2^2$$

which is a denoising objectif $\rightarrow$ analogous to denoising score matching

# Results

Score based and Diffusion models are now state-of-the-art in image generation

## Some ressources

https://lilianweng.github.io/posts/2021-07-11-diffusion-models/

https://yang-song.net/blog/2021/score/

https://dl.heeere.com/conditional-flow-matching/blog/
conditional-flow-matching/