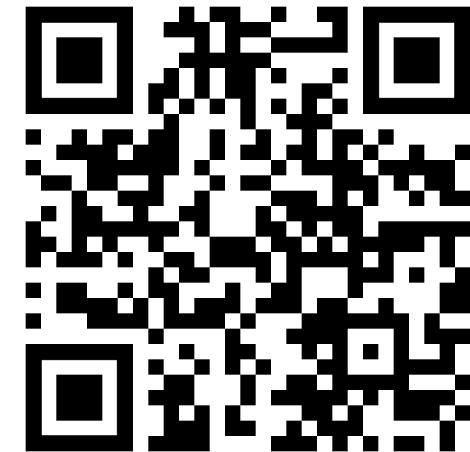


# Density Ratio Estimation with Conditional Probability Paths



Hanlin Yu<sup>1</sup>, Arto Klami<sup>1</sup>, Aapo Hyvärinen<sup>1</sup>, Anna Korba<sup>2</sup>, Omar Chehab<sup>2</sup>

1. University of Helsinki, Finland 2. ENSAE, CREST, IP Paris, France



## Problem statement

Given samples from two distributions,  $X_0 \sim p_0$  and  $X_1 \sim p_1$ , estimate the ratio  $\frac{p_1(\mathbf{x})}{p_0(\mathbf{x})}$ .

**Algorithm** [Choi et al., AISTATS 2022]

1. Interpolate samples:  $X_t = \sqrt{1-t^2}X_0(\mathbf{x}) + \sqrt{t^2}X_1(\mathbf{x})$ . The law  $p_t(\mathbf{x})$  is implicit.
2. Estimate the time score  $\partial_t \log p_t(\mathbf{x})$ .
3. Obtain the log ratio through numerical integration:  $\log \frac{p_1(\mathbf{x})}{p_0(\mathbf{x})} = \int_0^1 \partial_t \log p_t(\mathbf{x}) dt$

## Learning objectives for the time score

Original regression

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{x}, t)} [\lambda(t)(\partial_t \log p_t(\mathbf{x}) - s_{\boldsymbol{\theta}}(\mathbf{x}, t))^2] \quad \text{not explicit}$$

Integrate by parts  
TSM

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &= 2\mathbb{E}_{p_0(\mathbf{x})}[s_{\boldsymbol{\theta}}(\mathbf{x}, 0)] - 2\mathbb{E}_{p_1(\mathbf{x})}[s_{\boldsymbol{\theta}}(\mathbf{x}, 1)] \\ &+ \mathbb{E}_{p(t, \mathbf{x})}[2\partial_t s_{\boldsymbol{\theta}}(\mathbf{x}, t) + 2\dot{\lambda}(t)s_{\boldsymbol{\theta}}(\mathbf{x}, t) + \lambda(t)s_{\boldsymbol{\theta}}(\mathbf{x}, t)^2] \end{aligned} \quad \text{slow to differentiate}$$

Condition (ours)  
CTSM

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{x}, z, t)} [\lambda(t)(\partial_t \log p_t(\mathbf{x} | z) - s_{\boldsymbol{\theta}}(\mathbf{x}, t))^2] \quad \text{explicit}$$

Factorize (ours)  
CTSM-v

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{x}, z, t)} \left[ \lambda(t) \sum_{i=1}^D (\partial_t \log p_t(x^i | \mathbf{x}^{<i}, z) - s_{\boldsymbol{\theta}}^i(\mathbf{x}, t))^2 \right]$$

We also introduce the weighting function  $\lambda(t) \propto 1/|\partial_t \log p_t(\mathbf{x} | z)|$ .

**Theoretical guarantees** (modified): for  $K$  integration steps and  $N$  samples,

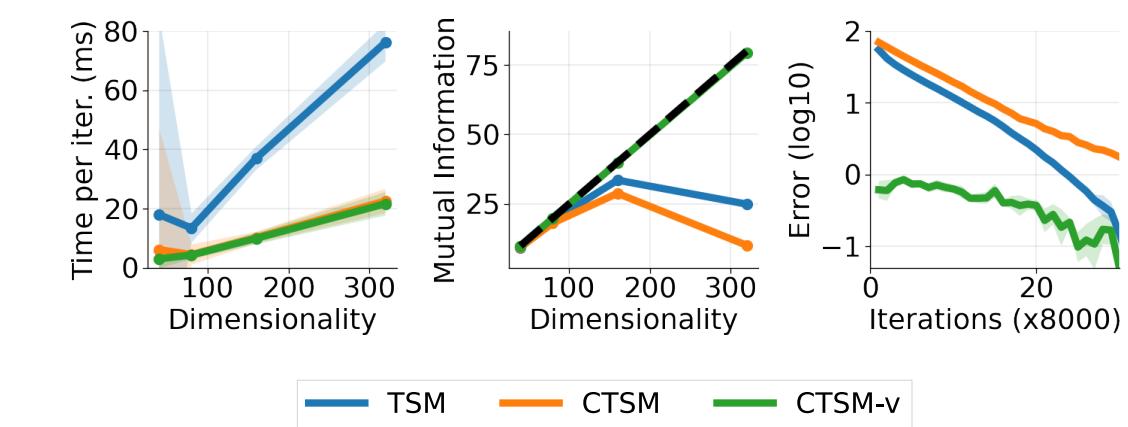
$$\mathbb{E}_{\hat{p}_1} \left\| \log \frac{p_1}{p_0} - \widehat{\log \frac{p_1}{p_0}} \right\|_{L^2(p_1)}^2 \leq \underbrace{\frac{1}{2K^2} \mathbb{E}_{p_1(\mathbf{x})}[L(\mathbf{x})^2]}_{\text{integral discretization error}} + \underbrace{\frac{2}{N} e(\theta^*, \lambda, p_t)}_{\text{score estimation error}} + o\left(\frac{1}{N}\right)$$

null if  $t \rightarrow \partial_t \log p_t(\mathbf{x})$  constant,  
i.e. Lipschitz constant  $L(\mathbf{x})$  is null

## Applications of density-ratio estimation

### Mutual information estimation

CTSM-v is faster and outperforms others especially in high dimensions.



**Likelihood estimation** (in bits per dimension, BPD). We use

$$\log p_1(\mathbf{x}) = \underbrace{\log p_0(\mathbf{x})}_{\text{Known}} + \underbrace{\int_0^1 \partial_t \log p_t(\mathbf{x}) dt}_{\text{Estimated}}$$

**Sample generation.** We convert the estimated time scores into space scores and plug them into popular score-based samplers.

$$\nabla \log p_t(\mathbf{x}) = \nabla \left( \underbrace{\log p_0(\mathbf{x})}_{\text{Known}} + \underbrace{\int_0^t \partial_s \log p_s(\mathbf{x}) ds}_{\text{Estimated}} \right)$$

Space	Methods	Approx. BPD	Time per step
Latent space	TSM	1.30	347 ms
	Ours	<b>1.26</b>	<b>58 ms</b>
Pixel space	TSM	unstable	1103 ms
	Ours	<b>1.03</b>	<b>142 ms</b>

Annealed Langevin sampler

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

Probability flow ODE sampler

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