

# Efficient Model Inference with Stein Density Ratio Estimation

Tractable Inference for Intractable Models

---

Song Liu (song.liu@bristol.ac.uk)<sup>1,4</sup>,

with Wittawat Jitkrittum<sup>2</sup>,

Takafumi Kanamori<sup>3</sup>,

and Yu Chen<sup>1</sup>

Special thanks to Carl Henrik Ek<sup>1</sup>

<sup>1</sup>University of Bristol, <sup>2</sup>Max-Planck Institute,

<sup>3</sup>Tokyo Institute of Technology, <sup>4</sup>Alan Turing Institute

# Paper/Code Available

- Paper: <https://arxiv.org/abs/1805.07454>
- Code: <https://github.com/lamfeeling/Stein-Density-Ratio-Estimation>

## 📖 Stein Density Ratio Estimation (SDRE) and Its Applications

### Reference:

*Song Liu, Takafumi Kanamori, Wittawat Jitkrittum, Yu Chen, Fisher Efficient Inference of Intractable Models, E-print: arXiv:1805.07454, To appear NeurIPS2019, 2019,*

### Install the `sdre` package

If you plan to modify our code (very likely, you will want to do so), it is best to install by:

1. Clone this repository
2. `cd` to the folder that you get, and install our package by (notice the dot at the end)

```
pip install -e .
```

# Table of contents

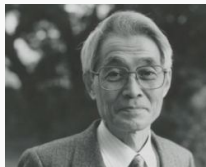
1. Intro
2. Stein Density Ratio Estimation
3. Intractable Model Inference
4. Conclusion

# Intro

---

# Problem: Model Inference

- Given **dataset**  $X_q := \{\mathbf{x}_q^{(i)}\}_{i=1} \sim q_x$  and **model**  $p(\mathbf{x}; \theta)$ ,
- Finding **model parameter**  $\theta$  that minimizes a statistical discrepancy (such as Kullback-Leibler (KL) divergence).
- Well studied problem: (**Wilk's** Thoerem; **Fisher's** Maximum Likelihood Estimation; **Akaike's** Information Criterion).



# Minimizing KL Divergence

## Minimize a KL divergence

1.  $\text{KL}[q_X|p_\theta] = \int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x};\theta)} d\mathbf{x} = - \int q(\mathbf{x}) \log p(\mathbf{x};\theta) d\mathbf{x} + C$
2.  $\int q(\mathbf{x}) \log p(\mathbf{x};\theta) d\mathbf{x} \approx \frac{1}{n} \sum_{i=1}^n \log p(\mathbf{x}_q^{(i)}; \theta).$
3.  $\hat{\theta} := \underset{\theta}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \log p(\mathbf{x}_q^{(i)}; \theta)$

## 3. is called **Maximum Likelihood Estimation (MLE)**

- **Normality:**  $\hat{\theta}$  follows a normal distribution.
- **Efficiency:**  $\hat{\theta}$  has the lowest asymptotic variance.

# Intractable Model

$$p(\mathbf{x}; \theta) := \frac{\bar{p}(\mathbf{x}; \theta)}{Z(\theta)}$$

- $\bar{p}(\mathbf{x}; \theta)$ : unnormalized density, which has a parametric form and can be easily calculated.
- $Z(\theta) := \int \bar{p}(\mathbf{x}; \theta) d\mathbf{x}$ : normalization term.
- $p(\mathbf{x}; \theta)$  may be **intractable**.
- $\bar{p}(\mathbf{x}; \theta)$  is so complicated that  $Z(\theta)$  **cannot be calculated**.
- We want to estimate  $\theta$ ; can we still minimize KL?

# Tractability of Kullback-Leibler Divergence

KL is computationally intractable for an intractable model:

$$\text{KL}[q_X \| p_\theta] := \mathbb{E}_q \left[ \log \frac{q(\mathbf{x})}{p(\mathbf{x}; \theta)} \right].$$

- $p(\mathbf{x}; \theta)$  cannot be evaluated.

## Question

Can we find a *tractable* surrogate of KL from  $p_\theta$  to  $q$ ?

- Fisher Divergence e.g. [Lyu, 2009, Hyvärinen, 2005]
- Kernel Stein Discrepancy (KSD) [Barp et al., 2019]

Can we find a surrogate that **still mimics behaviors of KL**?



# Approximating KL

**KL from  $p$  to  $q$  can be efficiently approximated**

- $\text{KL}[q||p]$  is an expectation over a log ratio  $q/p$ .
- If we could approximate the ratio  $q/p$ , we can use sample to approximate the expectation.

**Density Ratio Estimation:** Given two samples,

$$X_q := \{\mathbf{x}_q^{(i)}\}_{i=1}^{n_q} \sim q, X_p := \{\mathbf{x}_p^{(j)}\}_{j=1}^{n_p} \sim p,$$

Estimate  $\hat{r}$  as the ratio between  $q$  and  $p$  using  $X_q$  and  $X_p$ .

$$\text{KL}[q|p] = \mathbb{E}_q \left[ \log \frac{q(\mathbf{x})}{p(\mathbf{x})} \right] \approx \frac{1}{n_q} \sum_{i=1}^{n_q} \log \hat{r}(\mathbf{x}_q^{(i)}).$$

## Density Ratio Estimation (KLIEP) (1)

Let us introduce a parametric model for density ratio:  $r(\mathbf{x}; \delta)$ .

$\delta$  is fitted by minimizing a KL divergence  $\text{KL}[q\|r_\delta p]$ , i.e.,

$$\delta^* := \underset{\delta}{\operatorname{argmin}} \text{KL}[q\|r_\delta p] \quad \text{s.t.} \quad \int r(\mathbf{x}; \delta) p(\mathbf{x}) d\mathbf{x} = 1, \quad (1)$$

where  $\text{KL}[q\|r_\delta p] = \mathbb{E}_q[\log q(\mathbf{x})] - \mathbb{E}_q[\log r(\mathbf{x}; \delta)] - \mathbb{E}_q[\log p(\mathbf{x})]$

$$\approx -\frac{1}{n_q} \sum_{i=1}^{n_q} \log r(\mathbf{x}_q^{(i)}; \delta) + C,$$

and  $C$  is a constant [Sugiyama et al., 2008].

## Density Ratio Estimation (2)

We can also approximate the equality constraint in (1) using  $X_p$ :

$$\int r(\mathbf{x}; \delta) p(\mathbf{x}) d\mathbf{x} \approx \frac{1}{n_p} \sum_{j=1}^{n_p} r(\mathbf{x}_p^{(j)}; \delta).$$

Thus, (1) can be re-written as:

$$\begin{aligned} \hat{\delta} := \underset{\delta}{\operatorname{argmin}} & -\frac{1}{n_q} \sum_{i=1}^{n_q} \log r(\mathbf{x}_q^{(i)}; \delta) \\ \text{s.t.} & \frac{1}{n_p} \sum_{j=1}^{n_p} r(\mathbf{x}_p^{(j)}; \delta) = 1. \end{aligned} \tag{2}$$

Tractable, no need for MCMC sampling.

## KLIEP Approximated KL for Model Inference?

Given  $X_q$  and  $p(\mathbf{x}; \theta)$ , how to estimate  $q(\mathbf{x})/p_\theta(\mathbf{x})$ ?

We only have *one* sample  $X_q$  and a model  $p_\theta$ .

# Stein Density Ratio Estimation

---

# Stein Features

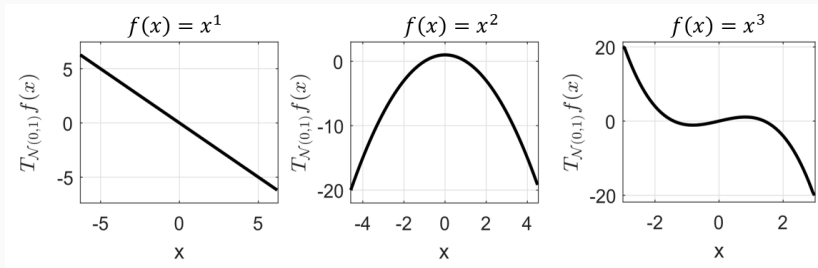
- Suppose  $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^b$  is a feature function.
  - e.g.  $\mathbf{f}(\mathbf{x}) = [1, x, x^2 \dots x^{b-1}]$ .
- Let  $T_\theta$  be a mapping:

$$T_\theta \mathbf{f}(\mathbf{x}) := [T_\theta f_1(\mathbf{x}), T_\theta f_2(\mathbf{x}), \dots, T_\theta f_b(\mathbf{x})],$$
$$T_\theta f_i(\mathbf{x}) := \langle \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}), \nabla_{\mathbf{x}} f_i(\mathbf{x}) \rangle + \text{tr} [\nabla_{\mathbf{x}}^2 f_i(\mathbf{x})],$$

where  $f_i$  is the  $i$ -th output of function  $\mathbf{f}$ .

- $T_\theta$  is called Stein Operator  
[Oates et al., 2017, Chwialkowski et al., 2016].
- Computing  $T_\theta \mathbf{f}(\mathbf{x})$  does **not** require the knowledge of  $Z(\theta)$ .

# Stein Features: Visualization



Stein feature with respect to  $p_\theta = \mathcal{N}(0, 1)$  defined on  $f_i(x) := x^i, i \in \{1, 2, 3\}$ .

- Stein Equality says, under **mild** regularity conditions:

$$\mathbb{E}_{p_\theta}[T_\theta \mathbf{f}(\mathbf{x})] = \mathbf{0}.$$

Proof see e.g., Lemma 5.1, [Chwialkowski et al., 2016]

- True for a large family of (continuously differentiable)  $\mathbf{f}$ ,
  - kernel smoother
  - neural network, with some activation functions.
  - polynomial,

w.r.t a wide range of  $p_\theta$ .



# Stein Density Ratio Model

We model density ratio function using Stein features:

$$r(\mathbf{x}; \delta) := \delta^\top T_\theta \mathbf{f}(\mathbf{x}) + 1$$

It can be seen that

$$\begin{aligned} \int p_\theta(\mathbf{x}) r(\mathbf{x}; \delta) d\mathbf{x} &= \mathbb{E}_{p_\theta} [\delta^\top T_\theta \mathbf{f}(\mathbf{x}) + 1] \\ &= \delta^\top \mathbb{E}_{p_\theta} [T_\theta \mathbf{f}(\mathbf{x})] + 1 = 1. \end{aligned}$$

The last equality is due to **Stein equality**.

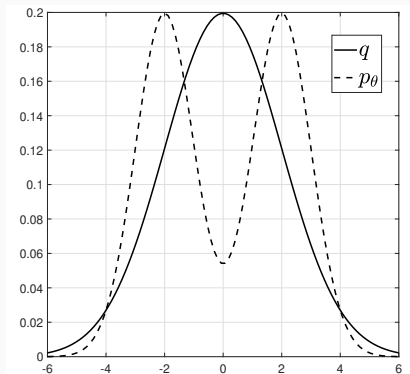
# Stein Density Ratio Estimation (SDRE), 1

Now we have a **Stein Density Ratio Estimator**.

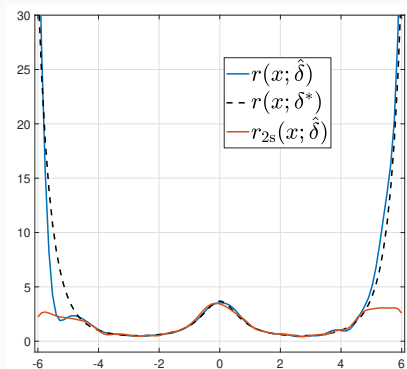
$$\begin{aligned}\hat{\delta} &:= \operatorname{argmax}_{\delta} \frac{1}{n_q} \sum_{i=1}^{n_q} \log r(\mathbf{x}; \delta) \\ &= \operatorname{argmax}_{\delta} \frac{1}{n_q} \sum_{i=1}^{n_q} \log \underbrace{\left[ \delta^\top T_{\theta} \mathbf{f}(\mathbf{x}_q^{(i)}) + 1 \right]}_{\ell(\delta; X_q, \theta)}\end{aligned}$$

- This is an unconstrained concave optimization.
- **Do not** need constraints, as it is automatically satisfied.

# One Sample vs. Two Sample Density Ratio Estimation



Density function  $q$  and  $p_\theta$



$r(x; \hat{\delta})$ ,  $r(x; \delta^*)$  and  $r_{2s}(x; \hat{\delta})$

# Intractable Model Inference

---

# Intractable Model Inference using “Discriminative-LE”

Ideally, MLE fits a  $\theta$  to  $X_q$  such that  $\text{KL}[q||p_\theta]$  is minimized.

- $\text{KL}[q||p_\theta]$  is **intractable** in many cases.
- Now we have a **tractable** approximation:

$$\text{KL}[q||p_\theta] \approx \mathbb{E}_q \left[ \log r(\mathbf{x}; \hat{\delta}) \right] \approx \max_{\delta} \ell(\delta; X_q, \theta)$$

- Naturally, we can do a “Discriminative-LE (DLE)”:

$$\hat{\theta} := \underset{\theta}{\operatorname{argmin}} \underbrace{\max_{\delta} \ell(\delta; X_q, \theta)}_{\text{discriminator}}$$

# Consistent?

- **Q:** Under what conditions, DLE is consistent?
- Specifically, we study the following estimator:

$$(\hat{\delta}, \hat{\theta}) := \arg \min_{\theta \in \Theta} \max_{\delta \in \Delta_{nq}} \ell(\delta, \theta).$$

- when  $q(\mathbf{x}) \equiv p(\mathbf{x}; \theta^*)$ ,
- when  $(\hat{\theta}, \hat{\delta})$  is a saddle point.

## Theorem (Consistency)

*Suppose  $p(\mathbf{x}; \theta^*) = q(\mathbf{x})$  and several regularity conditions hold.*

$$(\hat{\delta}, \hat{\theta}) \xrightarrow{\mathbb{P}} (\mathbf{0}, \theta^*).$$

MLE have many nice properties,  
such as computable confidence interval on parameters.  
However, computations is not possible when  $p_\theta$  is intractable.



# Asymptotic Normality of $\hat{\theta}$

$\hat{\theta}$  has a simple asymptotic distribution.

## Theorem (Asymptotic Normality of $\hat{\theta}$ )

*Suppose  $p(\mathbf{x}; \theta^*) = q(\mathbf{x})$  and several regularity conditions hold.*

$$\sqrt{n_q} \left( \theta^* - \hat{\theta} \right) \rightsquigarrow \mathcal{N} [0, \mathbf{V}],$$

where  $\mathbf{V} := - \left( \mathbb{E}_q [\mathbf{H}^*]_{\theta, \delta} \mathbb{E}_q [\mathbf{H}^*]_{\delta, \delta}^{-1} \mathbb{E}_q [\mathbf{H}^*]_{\delta, \theta} \right)^{-1}$ ,  
 $-\mathbb{E}_q [\mathbf{H}^*] := -\mathbb{E}_q [\nabla^2 \ell(\mathbf{0}, \theta^*)]$ .

# Asymptotic Efficiency (Asymptotic Variance)

## Theorem (Asymptotic Fisher Efficiency)

Given two Stein features  $T_{\theta^*}\mathbf{f}$ ,  $T_{\theta^*}\bar{\mathbf{f}}$ , if  $\text{span } T_{\theta^*}\mathbf{f} \subseteq \text{span } T_{\theta^*}\bar{\mathbf{f}}$ ,  $\mathbf{V}_{\bar{\mathbf{f}}} \preceq \mathbf{V}_{\mathbf{f}}$ . If  $\nabla_{\theta} \log p(\mathbf{x}; \theta^*) \in \text{span } T_{\theta^*}\mathbf{f}$ ,  $\mathbf{V}_{\mathbf{f}}$  attains a lowerbound

$$\mathbb{E}_q[\nabla_{\theta} \log p(\mathbf{x}; \theta^*) \nabla_{\theta} \log p(\mathbf{x}; \theta^*)^{\top}],$$

which is the Fisher information.

### Note:

The **smallest** asymptotic variance an unbiased estimator can reach [Cramér, 1946, Rao, 1945].

Other asymptotic theorems can be found in our paper.

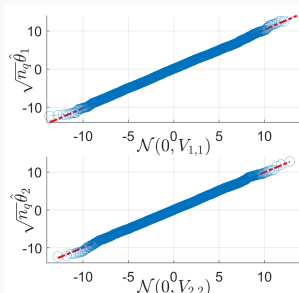
# Experiments

---

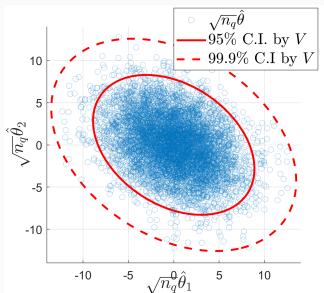
# Numerical Simulation: Distribution of $\hat{\theta}$

$$\bar{p}(\mathbf{x}; \theta) := \exp[\eta(\theta)^\top \psi(\mathbf{x})], \eta(\theta) := [-.5, .2, .6, 0, 0, 0, \theta]^\top,$$
$$\psi(\mathbf{x}) := [\sum_{i=1}^d x_i^2, \sum_{i=3}^d x_1 x_i, x_1 x_2, \tanh(\mathbf{x})]^\top, \mathbf{f}(\mathbf{x}) := [\tanh(\mathbf{x})]^\top.$$

$$q(\mathbf{x}) := p(\mathbf{x}; \mathbf{0})$$



(a) Marginal:  $\sqrt{n_q} \hat{\theta}$  vs. Asymptotic prediction



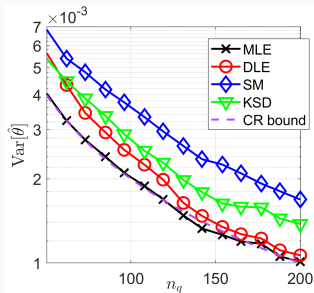
(b) Joint:  $\hat{\theta}$  vs. Asymptotic prediction

# Numerical Simulation: $\text{var} [\hat{\theta}]$

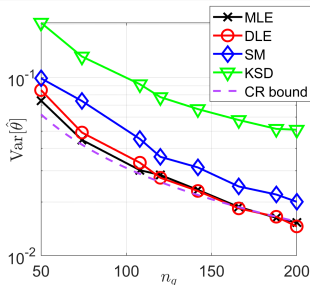
(Left): Gamma  $p(x; \theta) = \Gamma(5, \theta), \theta^* = 1$

(Right): Gaussian mixture model

$p(x; \theta) = .5\mathcal{N}(\theta, 1) + .5\mathcal{N}(1, 1), \theta^* = -1$



(c)  $\text{Var}(\hat{\theta})$ , Gamma dist.



(d)  $\text{Var}(\hat{\theta})$ , Gaussian mix.

Note: we use tractable models so Fisher info. can be easily calculated.

## Using Neural network as Sufficient Statistics

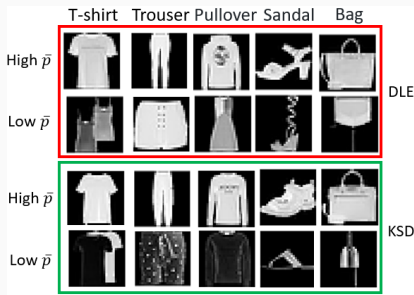
for  $i = 0 \dots 9$ ,  $\bar{p}(\mathbf{x}; \theta_i) := \exp[\theta_i^\top \psi(\mathbf{x})]$ ,  $\mathbf{x} \in \mathbb{R}^{784}$ ,  $\psi(\mathbf{x}) \in \mathbb{R}^{20}$  is a pre-trained 3-layer **DNN** on the full MNIST dataset with all digits.



Images with highest (upper) and lowest  $\log \bar{p}(\mathbf{x}; \hat{\theta})$  on each digit.

# Using Neural network as Sufficient Statistics

for  $i = 0 \dots 9$ ,  $\bar{p}(\mathbf{x}; \theta_i) := \exp[\theta_i^\top \psi(\mathbf{x})]$ ,  $\mathbf{x} \in \mathbb{R}^{784}$ ,  $\psi(\mathbf{x}) \in \mathbb{R}^{20}$  is a pre-trained 3-layer **DNN** on the entire Fashion MNIST dataset.



Images with highest (upper) and lowest  $\log \bar{p}(\mathbf{x}; \hat{\theta})$  on each digit.

## Conclusion




---



# Conclusion

- With the help of Stein Operator, we can estimate a density ratio between a parametric p.d.f and data distribution.
- Use  $\ell(\hat{\delta}, X_q, \theta)$  as a “tractable replacement” to  $\text{KL}[q|p_\theta]$  in model inference problems.
  - Discriminative Likelihood Estimation (DLE)
  - Asymptotic normality of  $\hat{\theta}$ .
  - Asymptotic variance of  $\hat{\theta}$  reaches Camera-Rao bound.
- Experiments on toy and real datasets show promising results.

Thank you very much!

-  Barp, A., Briol, F.-X., A.B., D., M., G., and L., M. (2019).  
**Minimum stein discrepancy estimators.**
-  Chwialkowski, K., Strathmann, H., and Gretton, A. (2016).  
**A kernel test of goodness of fit.**  
In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48, pages 2606–2615.
-  Cramér, H. (1946).  
**Mathematical methods of statistics.**  
Princeton university press.



Hyvärinen, A. (2005).

**Estimation of non-normalized statistical models by score matching.**

*Journal of Machine Learning Research*, 6:695–709.



Lyu, S. (2009).

**Interpretation and generalization of score matching.**

In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, pages 359–366. AUAI Press.



Oates, C. J., Girolami, M., and Chopin, N. (2017).

**Control functionals for monte carlo integration.**

*Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(3):695–718.



Rao, C. R. (1945).

**Information and the accuracy attainable in the estimation of statistical parameters.**

*Bulletin of the Calcutta Mathematical Society*, 37:81–91.



Sugiyama, M., Suzuki, T., Nakajima, S., Kashima, H., von Büna, P., and Kawanabe, M. (2008).

**Direct importance estimation for covariate shift adaptation.**

*Annals of the Institute of Statistical Mathematics*, 60(4):699–746.