

Bayesian Methods of Machine Learning.  
Project Presentation.  
Super-Samples from Kernel Herding

Georgii Novikov

October 2020

# Super-Samples from Kernel Herding

---

Weakly chaotic, non-linear dynamical system

$$\mathbf{x}_{t+1} = \arg \max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{w}_t, \phi(\mathbf{x}) \rangle \quad (1)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbb{E}_{\mathbf{x} \sim \mathbf{p}}[\phi(\mathbf{x})] - \phi(\mathbf{x}_{t+1})$$

$$\mathbf{x}_t \in \mathbb{R}^n, \phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Under some conditions, it is a greedy minimization of

$$\varepsilon_T^2 = \left\| \mu_{\mathbf{p}} - \frac{1}{T} \sum_{t=1}^T \phi(\mathbf{x}_t) \right\|^2, \text{ where } \mu_{\mathbf{p}} = \mathbb{E}_{\mathbf{x} \sim \mathbf{p}} \phi(\mathbf{x}) \quad (2)$$

# Theoretical Guarantees

---

## Theorem

$x_t$  in (1) is optimal on each step  $\Rightarrow$  error in (2) decreases at a rate  $\mathcal{O}(T^{-1})$ .

## Advantage

1. I.i.d samples have rate  $\mathcal{O}(T^{-\frac{1}{2}})$
2. MCMC converges even slower than  $\mathcal{O}(T^{-\frac{1}{2}})$  (due to positive correlations)

# Kernel Trick

---

We want to replace  $\phi(x)$  with  $k(x, x') = \langle \phi(x), \phi(x') \rangle$ :

1. Define  $w_0 = \mu := \mathbb{E}_{x \sim p}[\phi(x)]$
- 2.

$$\begin{aligned} x_{t+1} &= \arg \max_{x \in \mathcal{X}} \langle w_t, \phi(x) \rangle \\ &= \arg \max_{x \in \mathcal{X}} \langle w_0 + T \mathbb{E}_{x' \sim p}[\phi(x')] - \sum_{t=1}^T \phi(x_t), \phi(x) \rangle \\ &= \arg \max_{x \in \mathcal{X}} (T + 1) \mathbb{E}_{x' \sim p} k(x, x') - \sum_{t=1}^T k(x, x_t) \end{aligned}$$

## Interpretation

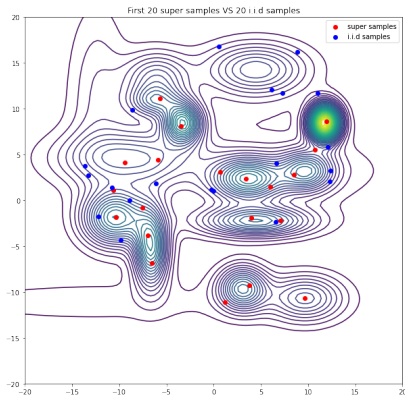
---

$$\begin{aligned}\varepsilon_T^2 &= \left\| \mu_p - \frac{1}{T} \sum_{t=1}^T \phi(\mathbf{x}_t) \right\|_{\mathcal{H}}^2 \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{x}' \sim p} k(\mathbf{x}, \mathbf{x}') - \frac{2}{T} \sum_{i=1}^T \mathbb{E}_{\mathbf{x} \sim p} k(\mathbf{x}, \mathbf{x}_t) + \frac{1}{T^2} \sum_{t, t'=1}^T k(\mathbf{x}, \mathbf{x}')\end{aligned}$$

This is a measure between distribution  $p$  and empirical distribution  $\hat{p}_T(\mathbf{x}) = \sum_{t=1}^T \delta(\mathbf{x}, \mathbf{x}_t)$ . And we still have a rate of convergence  $\mathcal{O}(T^{-1})$ !

# Experiment 0. Toy dataset

Mixture of 20 2D Gaussians:

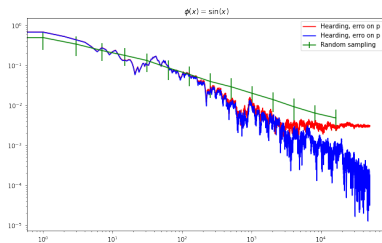
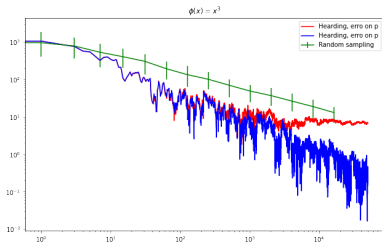
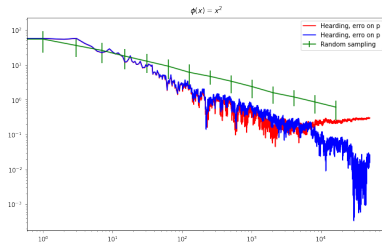
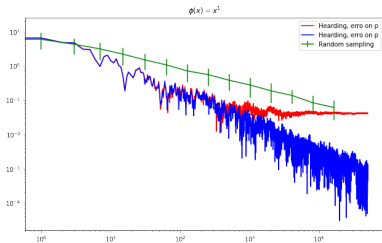


# Experiment 1. Empirical Matching

---

1.  $p(\mathbf{x})$  = mixture of 10 5D Gaussians
2.  $\mathcal{D} = 10^5$  i.i.d samples
3. Gaussian kernel  $k$  (with  $\sigma = 10$ )
4. Herding vs random samples on 4 functions of interest:
  - 4.1  $\phi(\mathbf{x}) = x^i, i \in 1, 2, 3$
  - 4.2  $\phi(\mathbf{x}) = \sin(\mathbf{x})$

# Experiment 1. Empirical Matching. Results





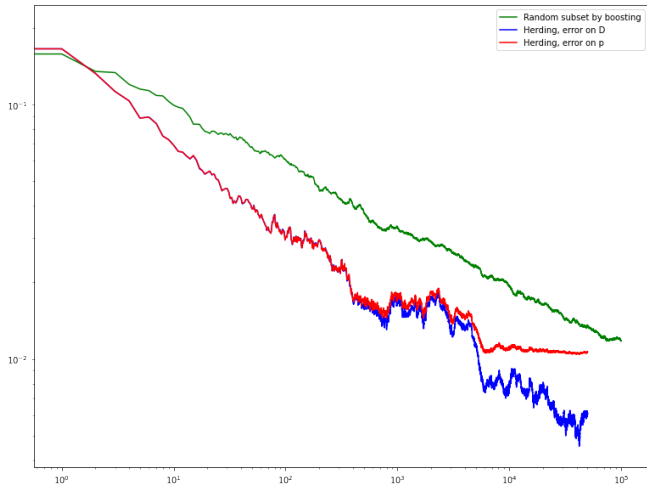
## Experiment 2. Approximating the Bayesian Posterior

---

- ▶ UCI spambase: 57 features. 4601 samples (3000 for train and 1601 for test).
- ▶ Whiten with PCA.
- ▶ Sample  $10^7$  logistic regression parameters from the poster distribution with gaussian prior with Metropolis-Hasting and subsample to  $10^5$  to reduce the autocorrelation.
- ▶ Whiten parameters with PCA.
- ▶ Herd super samples.
- ▶ Compare with metric

$$\text{RMSE}(S_T, D) = \frac{1}{N} \sum_{i=1}^N \left[ \frac{1}{T} \sum_{t=1}^T p(y_n | x_n, \theta_t) - \frac{1}{|D|} \sum_{d=1}^{|D|} p(y_n | x_n, \theta_d) \right]$$

## Experiment 2. Approximating the Bayesian Posterior. Result



# Problems

Dependence on the parameters of Gaussian kernel is not clear.

