

## INTRODUCTION

**Problem of interest:** How to better learn a complex and high-dimensional model in a big-data setting

### **Motivation:**

- Combine approaches and advantages from stochastic gradient MCMC (SG-MCMC) and stochastic optimization
- Stochastic optimization:
- computationally efficient iterations, fast convergence to a local optima
- Stochastic gradient MCMC:
- -computationally efficient iterations, slower convergence, able to explore the parameter space

Main idea: Begin with a preconditioned SG-MCMC algorithm, and gradually anneal the system temperature to zero such that it becomes a preconditioned stochastic optimization algorithm

**Advantages:** Our algorithm (Santa) has both an adaptive preconditioner and adaptive momentum, not available in existing algorithms.

**Software:** Code available at https://github.com/cchangyou/Santa.

## **STOCHASTIC OPTIMIZATION vs. SG-MCMC**

### **Stochastic optimization**

- Stochastic gradient descent (SGD) -basic stochastic optimization algorithm, without momentum and preconditioning
- SGD with momentum (SGD-M) -extending SGD with momentum
- RMSProp, ADAprop, Adadelta, . . . -extends SGD with a preconditioner
- Adam -extending SGD with both momentum and preconditioning

### **Stochastic gradient MCMC**

- Stochastic gradient Langevin dynamics (SGLD) -Sampling analog of SGD, without momentum and preconditioning
- Stochastic gradient Hamiltonian Monte Carlo (SGHMC)
- Sampling analog of SGD-M, with momentum
- Preconditioned stochastic gradient Langevin dynamics (pSGLD) -Sampling analog of RMSProp/Adagrad, with a preconditioner
- Multivariate stochastic gradient thermostats (mSGNHT) -Sampling with element-wise adaptive momentum, no obvious stochastic optimization analog

Table: SG-MCMC algorithms and stochastic optimization analags

Algorithms	SG-MCMC		Optimiza
Basic	SGLD	$\iff$	SGD
Precondition	pSGLD	$\iff$	RMSpro
Momentum	SGHMC	$\iff$	SGD-N
Thermostat	mSGNHT	$\iff$	Santa

# Bridging the Gap between Stochastic Gradient MCMC and Stochastic Optimization

Changyou Chen<sup>1</sup>, David Carlson<sup>2</sup>, Zhe Gan<sup>1</sup>, Chunyuan Li<sup>1</sup> and Lawrence Carin<sup>1</sup>  $^1$  Duke University, Durham, NC  $^2$  Columbia University, New York, NY

### ntion

\_\_\_\_\_

**Input:**  $\eta_t$  (learning rate),  $\sigma$ ,  $\lambda$ , burnin,  $\beta = \{\beta_1, \beta_2, \dots\} \rightarrow \infty$ ,  $\{\boldsymbol{\zeta}_t \in \mathbb{R}^p\} \sim N(\mathbf{0}, \mathbf{I}_p)$ . Initialize  $oldsymbol{ heta}_0$ ,  $oldsymbol{u}_0=\sqrt{\eta} imes N(0,I)$ ,  $oldsymbol{lpha}_0=\sqrt{\eta}C$ ,  $oldsymbol{v}_0=0$ . for t = 1, 2, ... do Evaluate  $\tilde{\boldsymbol{f}}_t \triangleq \nabla_{\boldsymbol{\theta}} \tilde{U}(\boldsymbol{\theta}_{t-1})$  on the  $t^{\text{th}}$  mini-batch.  $\boldsymbol{v}_t = \sigma \boldsymbol{v}_{t-1} + rac{1-\sigma}{N^2} \boldsymbol{f}_t \odot \boldsymbol{f}_t.$  $\boldsymbol{g}_t = 1 \oslash \sqrt{\lambda} + \sqrt{\boldsymbol{v}_t}.$ if t < burnin then  $\boldsymbol{\alpha}_t = \boldsymbol{\alpha}_{t-1} + (\boldsymbol{u}_{t-1} \odot \boldsymbol{u}_{t-1} - \eta/\beta_t).$  $\boldsymbol{u}_{t} = \frac{\eta}{\beta_{t}} \left( 1 - \boldsymbol{g}_{t-1} \oslash \boldsymbol{g}_{t} \right) \oslash \boldsymbol{u}_{t-1} + \sqrt{\frac{2\eta}{\beta_{t}}} \boldsymbol{g}_{t-1} \odot \boldsymbol{\zeta}_{t}.$ else  $oldsymbol{lpha}_t = oldsymbol{lpha}_{t-1}$  .  $u_t = 0.$ end if  $\boldsymbol{u}_t = \boldsymbol{u}_t + (1 - \boldsymbol{\alpha}_t) \odot \boldsymbol{u}_{t-1} - \eta \boldsymbol{g}_t \odot \boldsymbol{f}_t.$  $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + \boldsymbol{g}_t \odot \boldsymbol{u}_t.$ end for

 The Santa algorithm is based on the following stochastic differential equations, whose marginal distribution corresponds to the true posterior distribution of interest at temperature  $\frac{1}{\sqrt{2}}$ 

$$\begin{cases} d\boldsymbol{\theta} = G_1(\boldsymbol{\theta})\boldsymbol{p}dt \\ d\boldsymbol{p} = \left(-G_1(\boldsymbol{\theta})\nabla_{\boldsymbol{\theta}}U(\boldsymbol{\theta}) + G_1(\boldsymbol{\theta})(\boldsymbol{\Xi} - G_2(\boldsymbol{\theta}))\nabla_{\boldsymbol{\theta}}U(\boldsymbol{\theta}) + G_1(\boldsymbol{\theta})(\boldsymbol{\theta}) + G_1(\boldsymbol{\theta$$

- where  $m{Q} = \mathsf{diag}(m{p} \odot m{p})$ , w is standard Brownian motion,  $m{G}_1(m{ heta})$  and  $m{G}_2(m{ heta})$  are preconditioners.
- Santa algorithm is derived by solving (1) numerically with an increasing sequence of  $\beta$

## **Theorem (Convergence)**

With certain assumptions, the Santa algorithm converges in expectation to a global optima of a smooth function.

## EXPERIMENTS

### **Illustration**: Double-well potential example.



Figure: (Left) Double-well potential. (Right) The evolution of  $\theta$  using Santa and RMSprop algorithms.

## THE SANTA ALGORITHM

### Feedforward neural networks and convolutional neural networks

# exploration

*i* refinement

(1)

SGD Stoc NIN Maxou

Table:

## THEORY

$$- \mathbf{\Xi} \boldsymbol{p} + \frac{1}{\beta} \nabla_{\boldsymbol{\theta}} G_1(\boldsymbol{\theta}) \\ 7_{\boldsymbol{\theta}} G_2(\boldsymbol{\theta})) \, \mathrm{d}t + (\frac{2}{\beta} G_2(\boldsymbol{\theta}))^{\frac{1}{2}} \mathrm{d}w$$







This research was supported by ARO, DARPA, DOE, NGA, ONR and NSF.



Santa outperforms other algorithms in most cases.

e: Test error on N	<b>INIST</b> classifica	tion using FNN	I and CNN.
Algorithms	FNN-400	FNN-800	CNN
Santa	1.21%	1.16%	0.47%
Adam	1.53%	1.47%	0.59%
RMSprop	1.59%	1.43%	0.64%
SGD-M	1.66%	1.72%	0.77%
SGD	1.72%	1.47%	0.81%
SGLD	1.64%	1.41%	0.71%
BPB <sup>◊</sup>	1.32%	1.34%	
GD, Dropout <sup>¢</sup>	1.51%	1.33%	
toc. Pooling <sup>&gt;</sup>			0.47%
IN, Dropout $^{\circ}$			0.47%
xout, Dropout	*		0.45%

### **Recurrent neural networks**

Language modeling with an RNN; Tests on four publicly available datasets.

	0.00	0		
thms	Diano.	Nott.	Muse.	JSB.
nta	7.60	3.39	7.20	8.46
am	8.00	3.70	7.56	8.51
prop	7.70	3.48	7.22	8.52
)-M	8.32	3.60	7.69	8.59
D	11.13	5.26	10.08	10.81
	7.66	3.89	7.19	8.58
-M¢	8.37	4.46	8.13	8.71

Table: Test negative log-likelihood on 4 datasets.

### **GoogleNet for ImageNet classification**

## ACKNOWLEDGEMENTS