

# Bridging the Gap between Stochastic Gradient MCMC and Stochastic Optimization

Changyou Chen, David Carlson, Zhe Gan, Chunyuan Li,  
Lawrence Carin

May 2, 2016



# Outline

- 1 Introduction
- 2 The Santa algorithm
- 3 Experiments

# Introduction

- This paper is about how to better solve a complex, high-dimensional, nonlinear optimization problem in a big-data setting.
- Stochastic optimization:
  - computationally efficient, fast convergence, prone to local optimal
- Stochastic gradient MCMC:
  - computationally efficient, slower convergence, able to explore the parameter space
- Can we combine advantages from both?

# Stochastic optimization

- Stochastic gradient descent (SGD)
  - basic stochastic optimization algorithm, without considering neither momentum and preconditioning
- SGD with momentum (SGD-M)
  - extending SGD with momentum
- RMSProp, Adadelta ...
  - extending SGD with preconditioner
- Adam
  - extending SGD with both momentum and preconditioner

# Stochastic gradient MCMC

- Stochastic gradient Langevin dynamics (SGLD)
  - Bayesian analog of SGD, without considering neither momentum and preconditioning
- Stochastic gradient Hamiltonian Monte Carlo (SGHMC)
  - Bayesian analog of SGD-M, with momentum
- Preconditioned stochastic gradient Langevin dynamics (PSGLD)
  - Bayesian analog of RMSProp, with preconditioner
- Multivariate stochastic gradient thermostats (mSGNHT)
  - Bayesian sampling with adaptive momentum, does not have a stochastic optimization analog

# Bridging the gap

- We propose a stochastic optimization algorithm, **Santa**, that starts from a preconditioned version of mSGNHT, whose temperature is then annealed to zero.
- It has the advantages of both adaptive preconditioner and adaptive momentum.

**Table:** SG-MCMC algorithms and their optimization counterparts.

<b>Algorithms</b>	<b>SG-MCMC</b>		<b>Optimization</b>
<i>Basic</i>	SGLD	$\longleftrightarrow$	SGD
<i>Precondition</i>	pSGLD	$\longleftrightarrow$	RMSprop
<i>Momentum</i>	SGHMC	$\longleftrightarrow$	SGD-M
<i>Thermostat</i>	mSGNHT	$\longleftrightarrow$	Santa

# Outline

- 1 Introduction
- 2 The Santa algorithm**
- 3 Experiments

# The Santa algorithm

**Input:**  $\eta_t$  (learning rate),  $\sigma$ ,  $\lambda$ , *burnin*,  $\beta = \{\beta_1, \beta_2, \dots\} \rightarrow \infty$ ,

$$\{\zeta_t \in \mathbb{R}^p\} \sim N(\mathbf{0}, \mathbf{I}_p).$$

Initialize  $\theta_0$ ,  $\mathbf{u}_0 = \sqrt{\eta} \times N(0, I)$ ,  $\alpha_0 = \sqrt{\eta} C$ ,  $\mathbf{v}_0 = \mathbf{0}$  ;

**for**  $t = 1, 2, \dots$  **do**

Evaluate  $\tilde{\mathbf{f}}_t \triangleq \nabla_{\theta} \tilde{U}(\theta_{t-1})$  on the  $t^{\text{th}}$  mini-batch;

$$\mathbf{v}_t = \sigma \mathbf{v}_{t-1} + \frac{1-\sigma}{N^2} \tilde{\mathbf{f}}_t \odot \tilde{\mathbf{f}}_t ;$$

$$\mathbf{g}_t = \mathbf{1} \otimes \sqrt{\lambda + \sqrt{\mathbf{v}_t}} ;$$

**if**  $t < \textit{burnin}$  **then**

    /\* *exploration* \*/

$$\alpha_t = \alpha_{t-1} + (\mathbf{u}_{t-1} \odot \mathbf{u}_{t-1} - \eta / \beta_t);$$

$$\mathbf{u}_t = \frac{\eta}{\beta_t} (1 - \mathbf{g}_{t-1} \otimes \mathbf{g}_t) \otimes \mathbf{u}_{t-1} + \sqrt{\frac{2\eta}{\beta_t} \mathbf{g}_{t-1}} \odot \zeta_t$$

**else**

    /\* *refinement* \*/

$$\alpha_t = \alpha_{t-1}; \quad \mathbf{u}_t = \mathbf{0};$$

**end**

$$\mathbf{u}_t = \mathbf{u}_t + (1 - \alpha_t) \odot \mathbf{u}_{t-1} - \eta \mathbf{g}_t \odot \tilde{\mathbf{f}}_t; \quad \theta_t = \theta_{t-1} + \mathbf{g}_t \odot \mathbf{u}_t;$$

**end**



# Theory

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

$$\begin{cases} d\theta &= G_1(\theta)\mathbf{p}dt \\ d\mathbf{p} &= \left(-G_1(\theta)\nabla_{\theta}U(\theta) - \Xi\mathbf{p} + \frac{1}{\beta}\nabla_{\theta}G_1(\theta) \right. \\ &\quad \left.+ G_1(\theta)(\Xi - G_2(\theta))\nabla_{\theta}G_2(\theta)\right)dt + \left(\frac{2}{\beta}G_2(\theta)\right)^{\frac{1}{2}}d\mathbf{w} \\ d\Xi &= \left(\mathbf{Q} - \frac{1}{\beta}I\right)dt, \end{cases} \quad (1)$$

where  $\mathbf{Q} = \text{diag}(\mathbf{p} \odot \mathbf{p})$ ,  $w$  is standard Brownian motion,  $G_1(\theta)$  and  $G_2(\theta)$  are some preconditioners.

- Santa algorithm is derived by solving (1) numerically with an increasing sequence of  $\beta$ .

# Convergence properties

- The goal of Santa is to obtain  $\theta^*$  such that

$$\theta^* = \operatorname{argmin}_{\theta} U(\theta)$$

- $\{\theta_1, \dots, \theta_L\}$ : parameters collected from the algorithm.
- Sample average:  $\hat{U} \triangleq \frac{1}{L} \sum_{t=1}^L U(\theta_t)$ .
- Global optima:  $\bar{U} \triangleq U(\theta^*)$ .
- We study the convergence of the bias:  $|\mathbb{E}\hat{U} - \bar{U}|$ , and mean square error (MSE):  $\mathbb{E}(\hat{U} - \bar{U})^2$ .

# Convergence properties

## Theorem

*Under certain assumptions, the bias and MSE converge, for some constant  $C$  and  $D$ , and stepsize  $h$ , as*

$$\text{Bias} \leq C e^{-U(\theta^*)} \left( \frac{1}{L} \sum_{t=1}^L \int e^{-\beta_t \Delta U(\theta)} d\theta \right) + D \left( \frac{1}{Lh} + h^2 \right).$$

$$\text{MSE} \leq C^2 e^{-2U(\theta^*)} \left( \frac{1}{L} \sum_{t=1}^L \int e^{-\beta_t \Delta U(\theta)} d\theta \right)^2 + D^2 \left( \frac{1}{Lh} + h^4 \right).$$

- The first part characterizes the distance between the global optima and the annealing distributions  $e^{-\beta_t U(\theta)}$ ; the second part characterizes the distance between the sample average and the annealing posterior average. Both decrease with increasing  $L$ .

# Convergence properties

## Theorem

*Under certain assumptions, the bias and MSE converge, for some constant  $C$  and  $D$ , and stepsize  $h$ , as*

$$\text{Bias} \leq C e^{-U(\theta^*)} \left( \frac{1}{L} \sum_{t=1}^L \int e^{-\beta_t \Delta U(\theta)} d\theta \right) + D \left( \frac{1}{Lh} + h^2 \right).$$

$$\text{MSE} \leq C^2 e^{-2U(\theta^*)} \left( \frac{1}{L} \sum_{t=1}^L \int e^{-\beta_t \Delta U(\theta)} d\theta \right)^2 + D^2 \left( \frac{1}{Lh} + h^4 \right).$$

- The theorem indicates Santa converges in expectation closed to the global optima.

# Outline

- 1 Introduction
- 2 The Santa algorithm
- 3 Experiments**

# Illustration

- Optimizing the double-well potential:

$$U(\theta) = (\theta + 4)(\theta + 1)(\theta - 1)(\theta - 3)/14 + 0.5 .$$

- Start close to a local mode.
- RMSProp gets stuck, while Santa is able to jump out of the local mode.

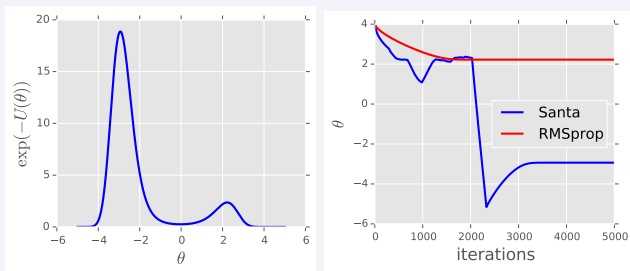


Figure: (Left) Double-well potential. (Right) The evolution of  $\theta$  using Santa and RMSprop algorithms.

# Feedforward neural networks and convolutional neural networks

- Detailed parameter setting is given in the paper.
- Santa outperforms other algorithms in most cases.

**Table:** Test error on MNIST classification using FNN and CNN.

Algorithms	FNN-400	FNN-800	CNN
Santa	<b>1.21%</b>	<b>1.16%</b>	<b>0.47%</b>
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%	—
SGD, Dropout <sup>◇</sup>	1.51%	1.33%	—
Stoc. Pooling <sup>▷</sup>	—	—	0.47%
NIN, Dropout <sup>◊</sup>	—	—	0.47%
Maxout, Dropout <sup>*</sup>	—	—	0.45%

# Recurrent neural networks (RNN)

- Language modeling with vanilla RNN.
- Test on four publicly available datasets.

**Table:** Test negative log-likelihood on 4 datasets.

Algorithms	Piano.	Nott.	Muse.	JSB.
Santa	<b>7.60</b>	<b>3.39</b>	<b>7.20</b>	<b>8.46</b>
Adam	8.00	3.70	7.56	8.51
RMSprop	7.70	3.48	7.22	8.52
SGD-M	8.32	3.60	7.69	8.59
SGD	11.13	5.26	10.08	10.81
HF <sup>◇</sup>	7.66	3.89	<b>7.19</b>	8.58
SGD-M <sup>◇</sup>	8.37	4.46	8.13	8.71



# GoogleNet for ImageNet classification

- These are preliminary results, did not report in the main text (included in the supplement).
- Use ILSVRC 2011 for training (ILSVRC 2012 has similar performance).
- Compared with SGD with momentum, other algorithms did not seem to work.
- Did not tune the parameters, use the default setting for GoogleNet provided in the Caffe package.
- Santa converges much faster than SGD-M.

# GoogleNet for ImageNet classification

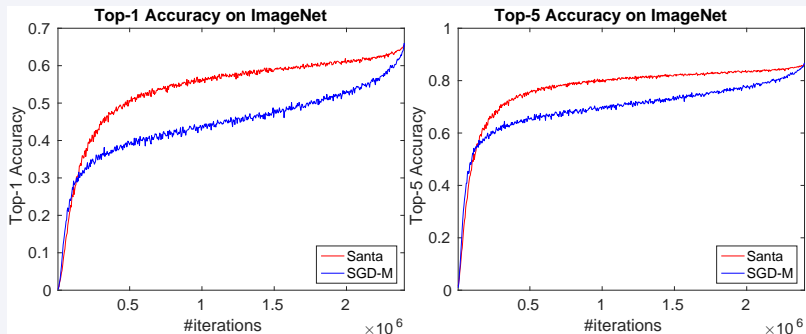


Figure: Santa vs. SGD with momentum on ImageNet.

# Code

- Code provided at <https://github.com/cchangyou/Santa>.
- Also provide a Caffe implementation.
- Welcome for feedbacks.

# Thanks for your attention

