

# Learning Weight Uncertainty with Stochastic Gradient MCMC for Shape Classification Chunyuan Li, Andrew Stevens, Changyou Chen, Yunchen Pu, Zhe Gan, Lawrence Carin

### Introduction

### Objective

- Weight uncertainty of deep neural networks (DNNs): posterior inference of weight distributions
- 2 Bring MCMC back to the community of computer vision to tackle "big visual/geometric data"
- Traditional MCMC: was popular in CV a decade ago, inclduing Gibbs sampling, HMC, MH, etc; but NOT scalable
- Scale up with Stochastic Gradient Markov chain Monte Carlo (SG-MCMC)

### Main Contributions

- We provide insights on the interpretation of Dropout from the perspective of SG-MCMC, which also allows the use of Batch-Normalization.
- 2 Applications to a wide range of shape classification problems demonstrate the advantages of SG-MCMC over optimization.



Figure: Illustration of Bayesian DNNs with a 2-layer model. All weights in Bayesian DNNs are represented as distributions using SG-MCMC (right figure); rather than having fixed values (left figure), as provided by classical stochastic optimization methods. The SG-MCMC learns correlated uncertainty jointly on all parameters, where (right) associated marginal distributions are depicted.

• Given data  $\mathcal{D} = \{ oldsymbol{d}_i \}_{i=1}^N$ ,  $oldsymbol{d}_i$  is *i.i.d.*; model parameters  $oldsymbol{ heta}$  $p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\boldsymbol{\theta}) \prod_{i=1}^{N} p(\boldsymbol{d}_i|\boldsymbol{\theta})$ 

$$\underbrace{P(\mathbf{o} \mid \mathbf{D})}_{\text{Posterior}} = \underbrace{P(\mathbf{o} \mid \mathbf{D})}_{\text{Prior}} = \underbrace{P(\mathbf{o} \mid \mathbf{D})}_{\text{Likelihood}}$$

For DNNs,  $d_i \triangleq (x_i, y_i)$ : input  $x_i \in \mathbb{R}^D$  and output  $y_i \in \mathcal{Y}$ . - Bayesian predictive estimate, for testing input  $ilde{m{x}}$ 

$$p(\tilde{y}|\tilde{\boldsymbol{x}}, \mathcal{D}) = \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})}[p(\tilde{y}|\tilde{\boldsymbol{x}}, \boldsymbol{\theta})] \approx \frac{1}{T} \Sigma_{t=1}^{T} p(\boldsymbol{\theta}|\mathcal{D})$$

In optimization,  $\theta_{MAP} = \operatorname{argmax} \log p(\theta | D)$ . The MAP approximates this expectation as

 $p(\tilde{y}|\tilde{\boldsymbol{x}}, \mathcal{D}) \approx p(\tilde{y}|\tilde{\boldsymbol{x}}, \boldsymbol{\theta}_{\mathsf{MAP}})$ 

Parameter uncertainty is ignored.

The predicted distribution of  $\tilde{y}$  may be viewed in terms of model averaging across parameters, based on the learned  $p(\theta | D)$ ; this should be contrasted with learning a single point estimate of  $\theta$  based on  $\mathcal{D}$ .

Duke University, Durham NC 27708, USA

### Algorithms

**Basic SG-MCMC algorithm** SGLD: Stochastic Gradient Langevin Dynamics [1]

### Algorithm 1: SGLD algorithm

Initialize:Random  $\theta_1$ ; for t = 1, 2, ..., T do %Estimate gradient from minibatch  $\mathcal{S}_{ extsf{t}}$  $\tilde{\boldsymbol{g}}_t \leftarrow \nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta}_t) + \frac{N}{|\mathcal{S}_t|} \sum_{i \in \mathcal{S}_t} \nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{d}_i | \boldsymbol{\theta}_t);$ %Parameter update  $\boldsymbol{\xi}_t \sim \mathcal{N}(0, \epsilon \mathbf{I});$  $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \frac{\epsilon}{2} \tilde{\boldsymbol{g}}_t + \boldsymbol{\xi}_t;$ end

where  $\epsilon$  is step-size, and  $S_t$  is the mini-batch

### More SG-MCMC algorithms

Table: SG-MCMC algorithms and their optimization counterparts. Algorithms in the same row share similar characteristics.

Algorithms	SG-MCMC	Optimization
Basic	SGLD	SGD
Preconditioning	pSGLD [2]	RMSprop/Adagrad/Adam
Momentum	SGHMC	momentum SGD
Thermostat	SGNHT	Santa [3]

# **Understanding Dropout**

### On the connection of SGLD and Dropout

For neural networks with the nonlinear function  $q(\cdot)$  and consecutive layers  $h_1$  and  $h_2$ , dropout and dropConnect:  $\boldsymbol{h}_2 = \boldsymbol{\xi}_0 \odot q(\boldsymbol{\theta} \boldsymbol{h}_1),$ Dropout: DropConnect:

where the injected noise  $\boldsymbol{\xi}_0$  can be binary valued with dropping rate p or its equivalent Gaussian : Binary noise:

Gaussian noise:

By combining dropConnect and Ga  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\xi}_0 \odot \boldsymbol{\theta}_t - \frac{\eta}{2} \tilde{\boldsymbol{f}}$ 

share the same form of update rule, with the only difference being that the level of injected noise is different.

Integration of SG-MCMC and Binary Dropout Accelerating SG-MCMC using Batch-Normalization

(1) $(\widetilde{y}|\widetilde{\boldsymbol{x}}, \boldsymbol{\theta}_t)$ 

(2)

- $\boldsymbol{h}_2 = q((\boldsymbol{\xi}_0 \odot \boldsymbol{\theta})\boldsymbol{h}_1),$

$$\begin{split} \boldsymbol{\xi}_{0} &\sim \operatorname{Ber}(p), \\ \boldsymbol{\xi}_{0} &\sim \mathcal{N}(1, \frac{p}{1-p}). \\ \text{aussian noise} \\ \boldsymbol{\tilde{f}}_{t} &= \boldsymbol{\theta}_{t} - \frac{\eta}{2} \tilde{\boldsymbol{f}}_{t} + \boldsymbol{\xi}_{0}' , \end{split} \tag{3}$$

where  $\boldsymbol{\xi}'_0 \sim \mathcal{N}\left(0, \frac{p}{(1-p)} \text{diag}(\boldsymbol{\theta}_t^2)\right)$ . Dropout/dropConnect and SGLD

## **Experiments: Shape Classification**

### . Datasets

	FNN	CNN		
<b>2D</b>	MNIST, Animal, 20 Newsgroups	MNIST, Caltech, Cifar10		
<b>3D</b>	Body Shape, Textured Shape	ModelNet		

Hand-crafted features as input of FNN, while raw data for CNN

### **II. Results and Observations** A thorough comparison

Accuracy of FNN on MNIST using a two-layer network (X-X) with ReLU. Please refer the paper for a lot more results.

Me

pSG SGL

RMS

SGE

pSG

SGL RMS

SGE

RMS

BPB BPE

# Very Deep Neural Networks

- The use of SG-MCMC or Dropout slows down learning initially. This is likely due  $\odot$ to the higher uncertainty imposed during learning, resulting in more exploration of the parameter space. Increased uncertainty, however, prevents
- overfitting and eventually results in improved performance.
- Exensive empirical results on adding gradient noise are also shown in [4].

### References

[1] Welling et al. Bayesian Learning via Stochastic Gradient Langevin Dynamics, ICML 2011 [2] Li et al, Preconditioned Stochastic Gradient Langevin Dynamics for Deep Neural Networks, AAAI 2016 [3] Chen et al, Bridging the Gap between Stochastic Gradient MCMC and Stochastic Optimization, AISTATS 2016 [4] Neelakantan et al, Adding Gradient Noise Improves Learning for Very Deep Networks, ICLR workshop 2016



ethods	Test	Error	(%)	
LD + Dropout	1.36	1.26	1.15	
D + Dropout	1.45	1.25	1.18	
Sprop + Dropout	1.35	1.28	1.24	
) + Dropout	1.51	1.33	1.36	
LD	1.45	1.32	1.24	
D	1.64	1.41	1.40	
Sprop	1.79	1.43	1.39	
)	1.72	1.47	1.47	
Sspectral	1.65	1.56	1.46	
3, Gaussian	1.82	1.99	2.04	
3, Scale mixture	1.32	1.34	1.32	

Networks 400-400 800-800 1200-1200

The testing error for the SG-MCMC methods are consistently lower than their corresponding stochastic optimization counterparts. This indicates that the weight uncertainty learned via SG-MCMC can improve performance.

Both of Dropout and SG-MCMC show their ability to regularize learning. By integrating SG-MCMC with Dropout, we obtain lower error.

### D: Dropout; BN: Batch Normalization



### Acknowledgements