

Jeremy Bernstein Caltech

Compressed optimisation for non-convex problems



Kamyar Azizzadenesheli UCI



signSGD

Yu-Xiang Wang UCSB/Amazon



Anima Anandkumar Caltech/Amazon







Snap gradient components to ±1

<u>Compressed</u> optimisation for non-convex problems

signSGD

Realistic for deep learning

STRUCTURE

Why care about signSGD?

Theoretical convergence results

Empirical characterisation of neural net landscape

Imagenet results





GRADIENT COMPRESSION WHY CARE?

. . . .



With 1/2 data



With 1/2 data





DISTRIBUTED SGD



With 1/2 data

With 1/2 data

SIGNSGD WITH MAJORITY VOTE

. . . .



With 1/2 data

With 1/2 data

.

COMPRESSION SAVINGS OF MAJORITY VOTE



Majority vote

SIGNSGD IS A SPECIAL CASE OF ADAM

signSGD sign $(g_k) = \frac{g_k}{\sqrt{g_k^2}}$

Signum sign $(g_k + \beta g_{k-1} + \beta^2 g_{k-2} + ...)$ (Sign momentum)

Adam
$$\frac{g_k + \beta g_{k-1} + \beta^2 g_{k-2} + \dots}{\sqrt{g_k^2 + \beta g_{k-1}^2 + \beta^2 g_{k-2}^2 + \dots}}$$

ADAM .WHY CARE?







UNIFYING ADAPTIVE GRADIENT METHODS + COMPRESSION

Sign descent

- weak theoretical foundation Ş
- incredibly popular (e.g. Adam)



Compressed descent

- weak theoretical foundation
- take pains to correct bias Ş
- empirically successful

Need to theory

Sign-based gradient compression?

STRUCTURE

Why care about signSGD?

Theoretical convergence results

Empirical characterisation of neural net landscape

Imagenet results

RE GD?





DOES SIGNSGD EVEN CONVERGE?

What might we fear?

- Might not converge at all
- Might have horrible dimension dependence
- Majority vote may give no speedup by adding extra machines

Our results

- ► It does converge
- Suggest these functions are typical in deep learning

Compression can be a free lunch

> We characterise functions where signSGD & majority vote are as nice as SGD

SINGLE WORKER RESULTS

Assumptions

- > Objective function lower bound f_*
- \succ Coordinate-wise variance bound $\overline{\sigma}$
- ► Coordinate-wise gradient Lipschitz **I**,

SGD gets rate
$$\mathbb{E}\left[\frac{1}{K}\sum_{k=0}^{K-1}\|g_k\|_2^2\right]$$
signSGD gets rate $\mathbb{E}\left[\frac{1}{K}\sum_{k=0}^{K-1}\|g_k\|_1\right]$

Define

- Number of iterations
- Number of backpropagations N







SINGLE WORKER RESULTS

Assumptions

- > Objective function lower bound f_*
- > Coordinate-wise variance bound $\overline{\sigma}$
- ► Coordinate-wise gradient Lipschitz



Define

- Number of iterations K
- Number of backpropagations N





MULTI WORKER RESULTS with M workers

Assumptions

- > Objective function lower bound f_*
- \blacktriangleright Coordinate-wise variance bound $\overline{\sigma}$
- Coordinate-wise gradient Lipschitz []

SGD gets rate

if gradient noise is unimodal symmetric majority vote gets

$$\mathbb{E} \left[\frac{1}{K} \sum_{k=0}^{K-1} \|g_k\|_2^2 \right]$$
$$\mathbb{E} \left[\frac{1}{K} \sum_{k=0}^{K-1} \|g_k\|_1 \right]$$



 $\leq \frac{1}{\sqrt{N}} \left| 2 \| \overrightarrow{L} \|_{\infty} (f_0 - f_*) + \frac{\| \overrightarrow{\sigma} \|_2^2}{\sqrt{M}} \right|$ $\frac{1}{\sqrt{N}} \left[\sqrt{\|\vec{L}\|_1} \left(f_0 - f_* + \frac{1}{2} \right) + 2 \frac{\|\vec{\sigma}\|_1}{\sqrt{M}} \right]^2$





STRUCTURE

Why care about signSGD?

Theoretical convergence results

Empirical characterisation of neural net landscape

Imagenet results



CHARACTERISING THE DEEP LEARNING LANDSCAPE EMPIRICALLY

Natural measure of density

$$\phi(\vec{v}) = \frac{\|\vec{v}\|_1^2}{d\|\vec{v}\|_2^2}$$

=1 for fully dense v

 ≈ 0 for fully sparse v



signSGD cares about gradient density
majority vote cares about noise symmetry



For large enough mini-batch size,

reasonable by Central Limit Theorem.

STRUCTURE

Why care about signSGD?

Theoretical convergence results

Empirical characterisation of neural net landscape

Imagenet results



SIGNUM IS COMPETITIVE ON IMAGENET



Performance very similar to Adam May want to switch to SGD towards end?

DOES MAJORITY VOTE WORK?

Cifar-10, Resnet-18





Test Accuracy



Jiawei Zhao NUAA



on server **pull** sign(\tilde{g}_m) **from** each worker **push** sign $\left[\sum_{m=1}^{M} \operatorname{sign}(\tilde{g}_m)\right]$ **to** each worker on each worker $x_{k+1} \leftarrow x_k - \delta \operatorname{sign} \left| \sum_{m=1}^M \operatorname{sign}(\tilde{g}_m) \right|$

Poster tonight! 6.15—9 PM @ Hall B #72