

# General Research Notes

Zachary Ross

May 3, 2022

# Contents

<b>I Dated Findings</b>	<b>3</b>
February 2, 2022 . . . . .	4
February 9, 2022 . . . . .	4
February 23, 2022 . . . . .	4
<b>II My Work</b>	<b>6</b>
<b>III Algorithms and Theory</b>	<b>7</b>
1 Monte Carlo and Las Vegas Randomized Algorithms for Systems and Control*: An Introduction [12]	8
2 Approximation Algorithms for k-Anonymity [1]	9
3 Improved Approximation Algorithms for Maximum Cut and Sat- isfiability Problems Using Semidefinite Programming [3]	10
3.1 Analysis . . . . .	11
4 Numerical Optimization [9]	13
4.6 Quasi-Newton Methods . . . . .	13
5 Treewidth, partial $k$ -trees, and chordal graphs [5]	14
<b>IV Deep Learning</b>	<b>15</b>
6 Efficient Per-Example Gradient Computations [4]	16
6.1 Commentary . . . . .	17
7 Understanding gradient clipping in incremental gradient meth- ods [10]	18
8 Pruning neural networks without any data by iteratively conserv- ing synaptic flow [11]	22
8.1 Pruning Algorithms . . . . .	22
8.2 Synaptic Saliency . . . . .	23

<b>CONTENTS</b>	<b>2</b>
8.3 Algorithm . . . . .	23
8.4 Commentary . . . . .	24
<b>9 Learning Low-rank Deep Neural Networks via Singular Vector Orthogonality Regularization and Singular Value Sparsification [13]</b>	<b>25</b>
<b>10 Stabilizing Gradients for Deep Neural Networks via Efficient SVD Parameterization [14]</b>	<b>27</b>
10.1 SVD and Orthogonality with Householder Reflectors . . . . .	28
<b>11 Federated Learning: Strategies for Improving Communication Efficiency [6]</b>	<b>30</b>
11.1 Structured Updates . . . . .	31
11.1.1 Low Rank . . . . .	31
11.1.2 Random Mask . . . . .	31
11.2 Sketched Updates . . . . .	31
11.2.1 Subsampling . . . . .	31
11.2.2 Probabilistic Quantization . . . . .	31
<b>12 Scalable Adaptive Stochastic Optimization Using Random Projections [7]</b>	<b>33</b>
<b>V Finance</b>	<b>36</b>
<b>13 An Introduction to Quantitative Finance [2]</b>	<b>37</b>
13.1 Interest Rates . . . . .	37
13.2 Forward Contracts and Forward Prices . . . . .	40
<b>14 A Random Walk Down Wall Street [8]</b>	<b>42</b>
<b>VI Differential Privacy</b>	<b>43</b>

# Part I

# Dated Findings

## February 2, 2022

In the work presented to me, Dr. Lee primarily uses pruning for gradient masking rather than parameter masking, although most studies I've read seem to have done the opposite. I initially assumed that given the standard linear layer function with masking

$$\mathbf{z} = (\mathbf{M} \odot \mathbf{W})\mathbf{h} + \mathbf{b} \quad (1)$$

the effect of the mask would still result in non-zero values in the gradient of the matrix. This is falsified by the derivation

$$\frac{\partial L}{\partial W_{i,j}} = \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial W_{ij}} = M_{ij} \frac{\partial L}{\partial z_i} h_j \quad (2)$$

or in vectorized form we have that

$$\frac{\partial L}{\partial \mathbf{W}} = \mathbf{M} \odot \left( \frac{\partial L}{\partial \mathbf{z}} (\mathbf{h})^\top \right) \quad (3)$$

which implies that the use of a forward mask guarantees that of a backward. The method used instead in Dr. Lee's examples was  $\mathbf{z} = \mathbf{W}\mathbf{h} + \mathbf{b}$  with masking only applied during the gradient (i.e. that of Equation 3). So what needs to be tested is whether or not we can apply masking solely during the forward phase i.e. implement Equation 1 with builds its gradient without the mask.

## February 9, 2022

Thorough testing has disproved our previous hypothesis. Strangely enough, the gradient masking is much more effective than weight masking for some odd reason.

## February 23, 2022

In the following discussion, I use beats per minute (BPM) to refer to heart rate. I only do it this way because I've already written it out and don't feel like changing it.

I've been pondering a bit recently on a hypothesis I've had, that being the perception of time is relative to the individual based on several factors. At first I believed that it's connected to BPM, but this was disproved by my addiction to coffee. I've noticed this relativity even in my day to day, some days seem faster than others while certain activities seem to be slower. For example, under the influence of coffee I feel that my day goes by much faster, which leads me to believe the effects of caffeine are misunderstood (which I will go into later). After or during an intense workout, my perception of time seems to slow way the hell down. I notice this because I'll be listening to music, and maybe it has to do with my percussionist background with intense focus on steady timing, but the songs pass by much slower. This is incredibly curious.

On other days where I don't quite feel myself, I notice instead that the music tends to be much quicker, often times leading me to feel like I'm missing a lot.

For the coffee thing, maybe its affect on BPM does slow down perception of time, but caffeine could effectively and periodically cut out clips of time, making

us believe that time is moving faster. I believe this affects productivity because it allows us to move from task to task seemingly "quicker". I don't believe it fully alters our productivity; rather, it seems like it allows us to reach the reward much quicker and thereby reinforce the thought process that we are "getting shit done". This activity/reward behavior most likely becomes a self-fulfilling cycle as the momentum continues, motivating us to continue doing tasks.

I've also considered how we would test these hypotheses because it's rather elusive. Say we tried to get a set of people to listen to a metronome prior to an energy intensive activity and post said activity and asked them to compare the two for which one is faster. For one, how would we get reliable information? If we were just to pick random people from a population, they may not even be capable of discerning between two tempos. I suppose we could do a preliminary exam where we vet out those who could not tell the difference even between minor changes in tempo. We could pick out trained musicians, but then the study is limited to musicians technically. I guess we would just have to do the vetting procedure and only select the upper echelon of people which can tell extremely minute differences in tempo after a select period of time.

Thoughts I had while working out, lol.

## **Part II**

# **My Work**

# **Part III**

# **Algorithms and Theory**

# Chapter 1

## Monte Carlo and Las Vegas Randomized Algorithms for Systems and Control\*: An Introduction [12]

**Abstract.** *In this paper, we present an introduction to Monte Carlo and Las Vegas randomized algorithms for systems and control. Specific applications of these algorithms include stability analysis, Lyapunov functions, and distributed consensus problems.*

## Chapter 2

# Approximation Algorithms for k-Anonymity [1]

**Abstract.** *We consider the problem of releasing a table containing personal records, while ensuring individual privacy and maintaining data integrity to the extent possible. One of the techniques proposed in the literature is  $k$ -anonymization. A release is considered  $k$ -anonymous if the information corresponding to any individual in the release cannot be distinguished from that of at least  $k-1$  other individuals whose information also appears in the release. In order to achieve  $k$ -anonymization, some of the entries of the table are either suppressed or generalized (e.g. an Age value of 23 could be changed to the Age range 20–25). The goal is to lose as little information as possible while ensuring that the release is  $k$ -anonymous. This optimization problem is referred to as the  $k$ -Anonymity problem. We show that the  $k$ -Anonymity problem is NP-hard even when the attribute values are ternary and we are allowed only to suppress entries. On the positive side, we provide an  $O(k)$ -approximation algorithm for the problem. We also give improved positive results for the interesting cases with specific values of  $k$  — in particular, we give a 1.5-approximation algorithm for the special case of 2-Anonymity, and a 2-approximation algorithm for 3-Anonymity.*

## Chapter 3

# Improved Approximation Algorithms for Maximum Cut and Satisfiability Problems Using Semidefinite Programming [3]

The initial problem attempted to be solved by this paper is creating an approximation algorithm for MAX-CUT to improve on the previous guarantee of  $\frac{3}{4}$  for an undirected graph  $G = (V, E)$  where  $|V| = n$  with nonnegative weights  $W_{ij} = W_{ji}$  i.e.  $\mathbf{W}$  denotes a weighted adjacency matrix.

We construct a cut of  $G$  by assigning to each  $i \in V$  a value  $y_i \in \{-1, 1\}$  and construct  $S \subset G$  by taking  $S = \{i | y_i = 1\}$ . Define the value of a cut as

$$w(S, \bar{S}) = \frac{1}{2} \sum_{i < j} W_{ij}(1 - y_i y_j) \quad (3.1)$$

The MAX-CUT integer quadratic program is then originally defined as

$$\text{MAX CUT}(G) = \max_{S \subset V} w(S, \bar{S}), \quad (3.2)$$

i.e. the goal is to find values  $y_i$  for all  $i \in V$  that maximize this equation.

The proposed relaxation is to substitute each  $y_i$  in Equation 3.1 with some vector  $\mathbf{v}_i \in \mathbb{R}^m$  ( $m \leq n$ ) such that  $\|\mathbf{v}_i\|^2 = 1$  so that each  $\mathbf{v}_i$  belongs to the  $m$ -dimensional unit sphere  $S_m$ . Under this relaxation, we instead construct  $S \subset G$  by selecting some  $\mathbf{r} \in S_m$  randomly and uniformly and take  $S = \{i | \langle \mathbf{v}_i, \mathbf{r} \rangle \geq 0\}$  i.e.  $\mathbf{r}$  determines a hyperplane which separates vertices of  $G$ . Then Equation 3.1 becomes

$$w(S, \bar{S}) = \frac{1}{2} \sum_{i < j} W_{ij}(1 - \langle \mathbf{v}_i, \mathbf{v}_j \rangle). \quad (3.3)$$

*Remark.* For a real, symmetric matrix  $A$ , the following are equivalent:

- $A$  is positive semidefinite;
- $\lambda(A) \subset \mathbb{R}^+ \cup \{0\}$  where  $\lambda(X)$  is the set of  $X$ 's eigenvalues; and
- $A = B^\top B$  for some  $B \in \mathbb{R}^{m \times n}$  where  $m \leq n$ .

The goal is now to determine optimal values for the unit vectors  $\mathbf{v}_1, \dots, \mathbf{v}_n$  and the dimensionality  $m$  for the vector space they reside. We can take  $\mathbf{B} = [\mathbf{v}_1 \ \dots \ \mathbf{v}_n]$  so that  $\mathbf{Y} = \mathbf{B}^\top \mathbf{B}$  is the Gram matrix of  $\mathbf{B}$ . Instead of solving for  $\mathbf{B}$  right away, we can instead note that each  $Y_{ij} = \langle \mathbf{v}_i, \mathbf{v}_j \rangle$  and that our previous constraint is now just the constraint  $Y_{ii} = 1$ . We also have that  $\mathbf{Y}$  is symmetric and PSD by construction. Semidefinite programming can be used to maximize the equation

$$\frac{1}{2} \sum_{i < j} W_{ij}(1 - Y_{ij}). \quad (3.4)$$

Rather than deal with each  $v_i$  directly, we can just use an approximate semidefinite programming algorithm to find an optimal  $\mathbf{Y}$  in  $O\left(\sqrt{n}\left(\log\left(\sum_{i < j} W_{ij}\right) + \log 1/\epsilon\right)\right)$  iterations (Alizadeh's adaptation of Ye's interior-point algorithm) with each iteration taking  $O(n^3)$ , where  $\epsilon > 0$  denotes the acceptable distance from the optimal value  $Z_P^*$  (i.e.  $Z_P^* - \epsilon$ ). Once an optimal  $\mathbf{Y}$  has been found, an incomplete Cholesky decomposition can be used to obtain  $\mathbf{v}_1, \dots, \mathbf{v}_n \in S_m$  s.t.

$$\frac{1}{2} \sum_{i < j} W_{ij}(1 - \langle \mathbf{v}_i, \mathbf{v}_j \rangle) \geq Z_P^* - \epsilon. \quad (3.5)$$

To find  $m$ , they note that any *extreme solution*, one which cannot be expressed as the strict convex combination of other feasible solutions, has at most rank  $l$  where

$$l \leq \frac{\sqrt{8n+1}-1}{2} < \sqrt{2n}. \quad (3.6)$$

Any other solution has  $m \leq n$ .

### 3.1 Analysis

The analysis breaks apart into three cases:

1. the general case with nonnegative edge weights,
2. an tighter bound when the weight of the cut constitutes a large portion of the overall weight,
3. and extension to negative weights.

The general case constitutes defining the overall expected weight as in

$$\mathbb{E}[W] = \frac{1}{\pi} \sum_{i < j} W_{ij} \Pr[\operatorname{sgn}(\langle \mathbf{v}_i, \mathbf{r} \rangle) \neq \operatorname{sgn}(\langle \mathbf{v}_j, \mathbf{r} \rangle)]. \quad (3.7)$$

They then state the following lemma:

**Lemma 3.1.**

$$\Pr [\operatorname{sgn}(\langle \mathbf{v}_i, \mathbf{r} \rangle) \neq \operatorname{sgn}(\langle \mathbf{v}_j, \mathbf{r} \rangle)] = \frac{1}{\pi} \arccos(\langle \mathbf{v}_i, \mathbf{v}_j \rangle), \quad (3.8)$$

which follows from the fact that the probability that these two values have the same sign is equivalent to the dihedral angle between the vectors, i.e. the probability that a randomly selected hyperplane bisects the two divided by all the possible hyperplanes that bisect the unit sphere. This results in the theorem

**Theorem 3.1.**

$$\mathbb{E}[W] = \frac{1}{\pi} \sum_{i < j} W_{ij} \arccos(\langle \mathbf{v}_i, \mathbf{v}_j \rangle). \quad (3.9)$$

Let  $\theta = \arccos(\langle \mathbf{v}_i, \mathbf{v}_j \rangle)$  and define

$$\alpha = \min_{0 \leq \theta \leq \pi} \frac{2}{\pi} \frac{\theta}{1 - \cos \theta}. \quad (3.10)$$

They then provide the following lemmas

**Lemma 3.2.** For  $-1 \leq y \leq 1$ ,  $\arccos(y)/\pi \geq \alpha \cdot \frac{1}{2}(1 - y)$ ,

which follows from substituting  $y = \cos \theta$ , and

**Lemma 3.3.**  $\alpha > 0.87856$

which follows from finding a lower bound for Equation 3.10. They use these to show that

**Theorem 3.2.**

$$\mathbb{E}[W] \geq \alpha \frac{1}{2} \sum_{i < j} W_{ij} (1 - \langle \mathbf{v}_i, \mathbf{v}_j \rangle). \quad (3.11)$$

For the case where the cut constitutes a large portion of the overall weight, we define  $W_{\text{tot}} = \sum_{i < j} W_{ij}$ ,  $h(t) = \arccos(1 - 2t)/\pi$ , and  $\gamma = \arg \min_{0 < t \leq 1} h(t)/t \approx 0.84458$ . The analysis then provides the theorem:

**Theorem 3.3.** Let

$$A = \frac{1}{W_{\text{tot}}} \frac{1}{2} \sum_{i < j} W_{ij} (1 - \langle \mathbf{v}_i, \mathbf{v}_j \rangle). \quad (3.12)$$

If  $A \geq \gamma$ , then

$$\mathbb{E}[W] \geq \frac{h(A)}{A} \frac{1}{2} \sum_{i < j} W_{ij} (1 - \langle \mathbf{v}_i, \mathbf{v}_j \rangle), \quad (3.13)$$

i.e. as  $A$  varies between  $\gamma$  and 1,  $h(A)/A$  varies between  $\alpha$  and 1. This provides us with a performance guarantee of  $h(A)/A - \epsilon$ .

## **Chapter 4**

# **Numerical Optimization [9]**

### **4.6 Quasi-Newton Methods**

# Chapter 5

## Treewidth, partial $k$ -trees, and chordal graphs [5]

(Read the paper and insert information here.)

**Definition 5.1.** A *tree decomposition* of a graph  $G = (V, E)$  is a pair

$$(\{X_i \mid i \in I\}, T = (I, M)) \quad (5.1)$$

where  $\{X_i \mid i \in I\}$  is a collection of subsets  $V$  (also called *bags*), and  $T$  is a tree s.t.

- $\cup_{i \in I} X_i = V$  (combining the vertices of all bags forms  $V$ )
- $(u, v) \in E \implies \exists i \in I : u, v \in X_i$  (if there is an edge connecting two vertices, then those vertices share at least one bag)
- $\forall v \in V : \{i \in I \mid v \in X_i\}$  induces a connected subtree of  $T$  (all bags involving a vertex  $v$  form a connected subtree).

*Remark.* The last statement in Definition 5.1 can be replaced with

- $i, k, j \in I$  and  $j$  is on the path from  $i$  to  $k$  in  $T \implies X_i \cap X_k \subseteq X_j$ .

**Lemma 5.1.** Let  $(\{X_i \mid i \in I\}, T = (I, M))$  be a tree decomposition of  $G = (V, E)$ , and let  $K \subseteq V$  be a clique in  $G$ . Then  $\exists i \in I : K \subseteq X_i$ .

# Part IV

# Deep Learning

# Chapter 6

## Efficient Per-Example Gradient Computations [4]

An important note about this paper is that it is *not* a derivation of the per-example gradient itself. Rather, it “describes an efficient technique for computing the *norm* of the gradient” with respect to the loss function.

Let each layer of a neural network be defined by the standard transformations

$$\begin{aligned} \mathbf{z}^{(i)} &= \mathbf{W}^{(i)} \mathbf{h}^{(i-1)} \\ \mathbf{h}^{(i)} &= \phi^{(i)}(\mathbf{z}^{(i)}) \end{aligned} \quad (6.1)$$

with the bias assumed to be an extra row/column of  $\mathbf{W}$  and loss function  $\mathcal{L}$ . Let  $\mathcal{B} \subset \mathcal{D}$  be a minibatch of the input dataset  $\mathcal{D}$ , and let  $\mathcal{B}^{(j)}$  be the  $j$ th example in the minibatch with  $\mathcal{L}^{(j)}$  corresponding to the loss of  $\mathcal{B}^{(j)}$  and cost function  $C = \sum_{j=1}^n \mathcal{L}^{(j)}$ . The per example gradient norm is computed with respect both to the example and the layer, such that the gradient norm for layer  $i$  and example  $j$  is

$$\left\| \frac{d\mathcal{L}^{(j)}}{d\mathbf{W}^{(i)}} \right\|^2 = \sum_{k,l} \left( \frac{\partial \mathcal{L}^{(j)}}{\partial W_{k,l}^{(i)}} \right)^2, \quad (6.2)$$

which is just the frobenius norm. Note that this equation can be used to calculate the  $L_2$  norm of the parameter gradient for an example or the  $L_2$  norm of the gradient for an individual weight matrix by summing over  $j$  or  $i$ , respectively.

The proposed method makes use of  $\mathbf{H}^{(i)}$  whose  $j$ th row contains the activation layer  $\mathbf{h}^{(i)}$  corresponding to  $\mathcal{B}^{(j)}$  i.e.  $\mathbf{H}_j^{(i)} = \phi^{(i)}(\mathbf{W}^{(i)} \phi^{(i-1)}(\dots(\mathbf{W}^{(1)} \mathcal{B}^{(j)}))\dots)$ . Likewise, define  $\mathbf{Z}^{(i)}$  for each  $\mathbf{z}^{(i)}$  and compute  $\bar{\mathbf{Z}} = \nabla_{\mathbf{Z}} C$ , which can be computed in a single pass using standard backpropagation. This gives us

$$\left\| \frac{\partial \mathcal{L}^{(j)}}{\partial \mathbf{W}^{(i)}} \right\|^2 = \left( \sum_k (\bar{Z}_{j,k}^{(i)})^2 \right) \left( \sum_k (H_{j,k}^{(i-1)})^2 \right) = \left\| \bar{\mathbf{Z}}_j^{(i)} \right\|^2 \left\| \mathbf{H}_j^{(i-1)} \right\|^2. \quad (6.3)$$

## 6.1 Commentary

This method works well with gradient clipping due to the efficient computation of gradient norms. This calculation allows for simple re-scaling of the gradients used in backwards propagation.

# Chapter 7

## Understanding gradient clipping in incremental gradient methods [10]

This paper considers the class of problems dealing with the minimization of

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{where} \quad f(x) = \frac{1}{m} \sum_{i=1}^m f_i(x), \quad (7.1)$$

where  $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$  are differentiable, non-convex functions. As a side note, this differs from notation that I'm used to, in that  $x$  is considered to be the set of parameters while  $i$ , or more importantly  $f_i$ , is a function of the parameters using the  $i^{\text{th}}$  component of the dataset to measure the loss. That is, the function  $f_i(x)$  determines the loss accumulated when using  $x$  as the set of parameters on input  $i$ .  $f$  is said to be the *objective function* while each  $f_i$  is called a *component function*.

The incremental method for SGD is

$$x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k) \quad k = 0, 1, \dots \quad (7.2)$$

where  $x_k$  is the values of the parameters at iteration  $k$ ,  $\alpha_k$  is the learning rate, and  $f_{i_k}$  is a uniform randomly chosen component function  $i_k \in \{1, 2, \dots, m\}$  of the dataset. The incremental method for IGC is

$$x_{k,i} = x_{k,i-1} - \alpha_k \nabla f_i(x_{k,i-1}) \quad i = 1, \dots, m \quad k = 0, 1, \dots \quad (7.3)$$

where the recursion is defined by base case  $x_{0,0}$  and iterative step  $x_{k,0} = x_{k-1,m}$  and is more cyclical by nature.

The clipping function  $\mathcal{C} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  with clipping threshold parameter  $\eta > 0$  is defined by

$$\mathcal{C}(g; \eta) = \min \left\{ 1, \frac{\eta}{\|g\|} \right\} \cdot g. \quad (7.4)$$

Our iterative methods can take advantage of clipping by clipping the gradient so that SGD with clipping is

$$x_{k+1} = x_k - \alpha_k \mathcal{C}(\nabla f_{i_k}(x_k); \eta), \quad (7.5)$$

and IGC with clipping is

$$x_{k,i} = x_{k,i-1} - \alpha_k \mathcal{C}(\nabla f_i(x_{k,i-1}); \eta), \quad (7.6)$$

and is written iteratively (using the notation  $x_k = x_{k,0} = x_{k-1,m}$ ) as

$$x_{k+1} = x_k - \alpha_k \sum_{i=1}^m \mathcal{C}(\nabla f_i(x_{k,i-1}); \eta), \quad (7.7)$$

For the theoretical understanding of the clipping function's effects on the convergence of both SGD and IGC, we make the following assumptions:

**Bounded below** there exists an  $f^*$  such that for all  $x \in \mathbb{R}^n$ ,  $f(x) \geq f^*$ ;

**Relaxed smoothness** there exist constants  $L_0, L_1 \in \mathbb{R}^+$  such that  $\|\nabla^2 f(x)\| \leq L_0 + L_1 \|\nabla f(x)\|$ .

The second condition becomes Lipschitz smoothness when  $L_1 = 0$ . Let  $f_k = f(x_k)$  and assume for simplicity from here on out that SGD and IGC have a constant step size  $\alpha$  and that IGC's notation for  $x_{k,i}$  can be simplified for  $x_{k,0}$  to  $x_k$ . The paper shows that convergence for the algorithm can be ensured if there exists a constant  $C$  such that

$$\left\langle \sum_{i=1}^m \mathcal{C}(\nabla f_i(x_k); \eta), \nabla f_k \right\rangle \geq C \|\nabla f_k\|^2 \quad (7.8)$$

i.e. if the magnitude of the clipped gradient projected in the direction of the regular gradient is bounded below by the magnitude of the regular gradient. This value is simplified further as follows.

Let  $g_{ki} = \nabla f_i(x_k)$  and note the following properties of  $g_{ki}$

$$\begin{aligned} \beta_{ki} &= \frac{\langle g_{ki}, \nabla f_k \rangle}{\|\nabla f_k\|^2} = \frac{\|g_{ki}\|}{\|\nabla f_k\|} \cos \theta_{ki} \\ g_{ki} &= \beta_{ki} \nabla f_k + t_{ki} \quad \text{where} \end{aligned} \quad . \quad (7.9)$$

$$t_{ki} \in \nabla f_k^\perp = \{z | \langle z, \nabla f_k \rangle = 0\}$$

Since  $\frac{1}{m} \sum_i^m g_{ki} = \nabla f_k$ , it is a necessary condition that  $\sum_i^m \beta_{ki} = m$  and  $\sum_i^m t_{ki} = 0$ . Let  $\gamma_{ki} = \min\{1, \eta / \|g_{ki}\|\}$  and note that  $\mathcal{C}(g_{ki}; \eta) = \gamma_{ki} g_{ki}$ . The paper then shows that

$$\left\langle \sum_{i=1}^m \mathcal{C}(\nabla f_i(x_k); \eta), \nabla f_k \right\rangle = \sum_{i=1}^m \gamma_{ki} \beta_{ki} \|\nabla f_k\|^2. \quad (7.10)$$

so that Equation 7.8 becomes

$$\sum_{i=1}^m \gamma_{ki} \beta_{ki} \geq C, \quad (7.11)$$

that is, *the necessary condition for convergence is that the sum in Equation 7.11 is greater than some positive scalar C.*

There are then essentially two conditions, either of which will ensure the convergence of our algorithm:

1. all  $g_{ki}$  can be partitioned into pairs so that if  $g_{kj}$  and  $g_{kl}$  are a pair then  $\beta_{kj} + \beta_{kl} > 0$  and  $\|t_{kj}\| = \|t_{kl}\|$ ;
2. all  $\beta_{ki} > 0$  or equivalently all  $\cos \theta_{ki} > 0$ .

Under the second condition, the paper then shows that a lower bound can be derived making use of the minimum and maximum gradient magnitudes w.r.t. an example  $i$ . That is, if we let  $m_g = \min_i \|g_{ki}\|$  and  $M_g = \max_i \|g_{ki}\|$ , and we let constants  $c_1$  and  $c_2$  be defined by the inequalities

$$0 < c_1 \leq \frac{1}{m} \sum_{i=1}^m \cos \theta_{ki} \quad \text{and} \quad 0 < c_2 \leq \frac{m_g}{M_g}, \quad (7.12)$$

then the lower bounds for Equation 7.8 are then divided into three scenarios:

1. if  $\eta \geq M_g$ , then

$$\left\langle \sum_{i=1}^m \mathcal{C}(\nabla f_i(x_k); \eta), \nabla f_k \right\rangle = m \|\nabla f_k\|^2; \quad (7.13)$$

2. if  $\eta \leq m_g$ , then

$$\left\langle \sum_{i=1}^m \mathcal{C}(\nabla f_i(x_k); \eta), \nabla f_k \right\rangle = \left( \eta \sum_{i=1}^m \cos \theta_{ki} \right) \|\nabla f_k\|; \quad (7.14)$$

3. and if  $m_g < \eta < M_g$ , then

$$\left\langle \sum_{i=1}^m \mathcal{C}(\nabla f_i(x_k); \eta), \nabla f_k \right\rangle \geq m \min\{\eta c_1, c_2 \|\nabla f_k\|\} \|\nabla f_k\|. \quad (7.15)$$

The paper concludes by providing the following theorems on the behavior of SGD and IGC with gradient clipping.

**Theorem 7.1.** *Assuming that we have a lower bound and relaxed smoothness, if*

$$\alpha \leq \frac{c_1}{4\eta L_1} \quad (7.16)$$

*then the iterates  $\{x_k\}$  generated by Equation 7.5 satisfy*

$$\frac{1}{T+1} \sum_{k=0}^T \min\{\eta c_1, c_2 \|\nabla f_k\|\} \|\nabla f_k\| \leq \frac{2(f_0 - f^*)}{(T+1)\alpha} + 5\alpha\eta^2 L_0 + \frac{\alpha\eta^3 L_1 c_1}{c_2}. \quad (7.17)$$

**Theorem 7.2.** Assuming that we have a lower bound, relaxed smoothness, and each component function also satisfies relaxed smoothness i.e.

$$\|\nabla^2 f_i(x)\| \leq L_{0i} + L_{1i} \|\nabla f_i(x)\|, \quad i = 1, \dots, m; \quad (7.18)$$

if

$$\alpha \leq \min \left\{ \frac{c_1}{2m(2\eta L_1 + G)}, \frac{1}{5mL_{0i}}, \frac{1}{8m\eta L_1 i} \right\}, \quad i = 1, \dots, m; \quad (7.19)$$

then the iterates  $\{x_k\}$  generated by Equation 7.7 satisfy

$$\frac{1}{T+1} \sum_{k=0}^T \min\{\eta c_1, c_2 \|\nabla f_k\|\} \|\nabla f_k\| \leq \frac{2(f_0 - f^*)}{(T+1)\alpha m} + 5\alpha m \eta^2 L_0 + \frac{\alpha m \eta^2 c_1 (2\eta L_1 + G)}{2c_2}, \quad (7.20)$$

where

$$G = 5 \max_{1 \leq i \leq m} L_{0i} + 16\eta \max_{1 \leq i \leq m} L_{1i}. \quad (7.21)$$

# Chapter 8

## Pruning neural networks without any data by iteratively conserving synaptic flow [11]

This paper, while theoretically captivating, is lackluster with providing mathematical commentary on the concepts and proofs. It may be worthwhile to go back and either note the more mathematical definitions of the paper, or derive my own definitions and reprove their claims using a more solid theoretical foundation. This may even set a better groundwork for future analysis.

### 8.1 Pruning Algorithms

Pruning algorithms are generally defined by

1. scoring parameters by some metric, and
2. masking parameters according to their scores.

The latter is generally done by removing the parameters (e.g. making the value zero via Hadamard product). This can either be done via global masking or layer-masking, with global-masking performing better but suffering from *layer-collapse*, which is when an entire layer is masked.

Let  $\theta$  be a collection of network parameters and  $\theta_{\text{prune}}$  be the parameters remaining after pruning. The *compression ratio*  $\rho$  of a pruning algorithm is

$$\rho = \frac{|\theta|}{|\theta_{\text{prune}}|} \tag{8.1}$$

. We define  $\rho_{\max}$  as the maximal possible compression ratio for a *network* that doesn't lead to layer collapse and the  $\rho_{\text{cr}}$  as the maximal compression ratio a given

*algorithm* can achieve without inducing layer collapse. Note that the distinction in the two is  $\rho_{\max}$  is maximal for a network while  $\rho_{\text{cr}}$  is maximal for an algorithm. These definitions motivate the following axiom:

**Axiom 8.1** (Maximal Critical Compression). *For any pruning algorithm and any network, we should always have  $\rho_{\text{cr}} = \rho_{\max}$ .*

## 8.2 Synaptic Saliency

*Synaptic saliency* is a class of score metrics defined by

$$\mathcal{S}(\theta) = \frac{\partial \mathcal{R}}{\partial \theta} \odot \theta \quad (8.2)$$

where  $\mathcal{R}$  is a scalar loss function of the output  $y$  of a feed-forward network parameterized by  $\theta$ . This metric satisfies two conservation laws:

**Theorem 8.1** (Neuron-wise Conservation of Synaptic Saliency). *For a feedforward neural network with continuous, homogeneous activation functions  $\phi(x) = \phi'(x)x$ , let  $j$  be the index of a hidden neuron in layer  $i$ . The sum of the synaptic saliency for the incoming parameters to a hidden neuron  $\left( \mathcal{S}_j^{(l)\text{in}} = \left\langle \frac{\partial \mathcal{R}}{\partial \theta_j^{(l)}}, \theta_j^{(l)} \right\rangle \right)$  is equal to the sum of the synaptic saliency for the outgoing parameters from the hidden neuron  $\left( \mathcal{S}_j^{(l)\text{out}} = \left\langle \frac{\partial \mathcal{R}}{\partial \theta_{:,j}^{(l+1)}}, \theta_{:,j}^{(l+1)} \right\rangle \right)$ .*

**Theorem 8.2** (Network-wise Conservation of Synaptic Saliency). *The sum of the synaptic saliency across any set of parameters that exactly separates the input neurons  $x$  from the output neurons  $y$  of a feedforward neural network with homogeneous activation functions equals  $\left\langle \frac{\partial \mathcal{R}}{\partial x}, x \right\rangle = \left\langle \frac{\partial \mathcal{R}}{\partial y}, y \right\rangle$ .*

## 8.3 Algorithm

For the following theorem, we define the *prune size* to be the total score for the parameters pruned at any iteration and the *cut size* to be the total score for an entire layer.

**Theorem 8.3.** *If a pruning algorithm with global-masking assigns positive scores that respect layer-wise conservation, and if the prune size is strictly less than the cut size whenever possible, then the algorithm satisfies the Maximal Critical Compression axiom.*

The pruning algorithm defined in this paper uses a loss function

$$\mathcal{R}_{\text{SF}} = \mathbf{1}^T \left( \prod_{l=1}^L |\theta^{(l)}| \right) \mathbf{1} \quad (8.3)$$

so that for a fully connected network ( $f_\theta(x) = \theta^{(L)} \dots \theta^{(1)} \mathbf{x}$ ) the Synaptic Flow score for a parameter  $\theta_{i,j}^{(l)}$  is

$$\mathcal{S}_{\text{SF}}(\theta_{i,j}^{(l)}) = \left[ \mathbf{1}^\top \prod_{k=l+1}^L |\theta^{(k)}| \right]_i |\theta_{i,j}^{(l)}| \left[ \prod_{k=1}^{l-1} |\theta^{(k)}| \mathbf{1} \right]_j \quad (8.4)$$

i.e. the portion of the  $l_1$ -path norm the network has through this parameter. This contributes to the development of Algorithm 1.

---

**Algorithm 1** Iterative Synaptic Flow Pruning (SynFlow)

---

**Require:** network  $f_\theta$ , compression ratio  $\rho$ , iteration steps  $n$

```

 $\mu = \mathbf{1}$                                      ▷ Initialize binary mask
for  $k$  in  $[1, \dots, n]$  do
     $\theta_\mu \leftarrow \mu \odot \theta_0$            ▷ Mask parameters
     $\mathcal{R} \leftarrow \mathbf{1}^\top \left( \prod_{l=1}^L |\theta_\mu^{(l)}| \right) \mathbf{1}$    ▷ Evaluate SynFlow objective
     $\mathcal{S} = \frac{\partial \mathcal{R}}{\partial \theta_\mu} \odot \theta_\mu$           ▷ Compute SynFlow score
     $\tau \leftarrow (1 - \rho^{-k/n})$  percentile of  $\mathcal{S}$       ▷ Find threshold
     $\mu \leftarrow (\tau < \mathcal{S})$                          ▷ Update mask
end for
return  $f_{\mu \odot \theta}$                            ▷ Return masked network

```

---

## 8.4 Commentary

A thought I had: the loss function defined in Equation 8.3 seems to almost indicate a just-past-the-minimum metric for the predefined network. It effectively is just the sum of the entries in the product of the network space. Although I agree with the data agnostic approach, I feel there should be a more suitable loss function which should seek to preserve this score rather than minimize it, or at the very least, should compare this score against some optimum, i.e. the optimum score for that network under those parameters. I also feel that the prune size and cut size should have a greater impact over the outcome of the resulting network, such that the compression ratio is fairly equal for all layers.

**Research Question 8.1.** Does SynFlow, or other pruning methods for that matter, produce the same network given different initial parameterizations? Furthermore, might we be able to use a genetic algorithm to find the “lottery ticket”?

One thing that needs to be considered with this question is that it’s entirely possible that isomorphic networks could be produced, so there may need to be some kind of check that occurs post pruning that confirms whether or not this network has been discovered already.

## Chapter 9

# Learning Low-rank Deep Neural Networks via Singular Vector Orthogonality Regularization and Singular Value Sparsification [13]

This paper focuses on an improvement on Deep Neural Networks (specifically ones involving convolutional layers) by utilizing the SVD of a convolution and sparsifying the underlying singular values. This is done by representing the kernel  $\mathbf{K}$  of a convolution as a 4-D tensor with dimension  $\mathbf{K} \in \mathbb{R}^{n \times c \times w \times h}$  where  $n$  is the number of filters (output channels),  $c$  is the number of input channels, and  $w, h$  are the width and height of the kernel respectively. They break the kernel down into the *channel-wise decomposition* and the *spatial-wise decomposition*.

Under channel-wise decomposition,  $\mathbf{K}$  is reshaped into a 2-D matrix  $\hat{\mathbf{K}} \in \mathbb{R}^{n \times cwh}$ , which can be decomposed with SVD into  $\mathbf{U} \in \mathbb{R}^{n \times r}$ ,  $\mathbf{V} \in \mathbb{R}^{cwh \times r}$ , and  $\mathbf{s} \in \mathbb{R}^r$  where  $\mathbf{U}, \mathbf{V}$  are unitary and  $r = \min(n, cwh)$ . The convolution then becomes the multiplication of convolutions  $\mathbf{K}_1 \in \mathbb{R}^{r \times c \times w \times h}$  (reshaped from  $\text{diag}(\sqrt{\mathbf{s}}) \mathbf{V}^\top$ ) and  $\mathbf{K}_2 \in \mathbb{R}^{c \times r \times 1 \times 1}$  (reshaped from  $\mathbf{U} \text{diag}(\sqrt{\mathbf{s}})$ ), a simplification discovered by [15].

Under spatial-wise decomposition,  $\mathbf{K}$  is reshaped into a 2-D matrix  $\hat{\mathbf{K}} \in \mathbb{R}^{nw \times ch}$ , which can be decomposed with SVD into  $\mathbf{U} \in \mathbb{R}^{nw \times r}$ ,  $\mathbf{V} \in \mathbb{R}^{ch \times r}$ , and  $\mathbf{s} \in \mathbb{R}^r$  where  $\mathbf{U}, \mathbf{V}$  are unitary and  $r = \min(n, cwh)$ . The convolution then becomes the multiplication of convolutions  $\mathbf{K}_1 \in \mathbb{R}^{r \times c \times 1 \times h}$  (reshaped from  $\text{diag}(\sqrt{\mathbf{s}}) \mathbf{V}^\top$ ) and  $\mathbf{K}_2 \in \mathbb{R}^{c \times r \times 1 \times 1}$  (reshaped from  $\mathbf{U} \text{diag}(\sqrt{\mathbf{s}})$ ).

In either case, the SVD decomposition of the matrix is stored in memory rather

Figure 9.1: Effect of applying different sparsity-inducing regularizers during SVD training. All models are achieved with Spatial-wise decomposition.

than the convolution, and the convolution is reconstructed from these matrices.<sup>o</sup>

Since the assurance that  $\mathbf{U}, \mathbf{V}$  are orthogonal is a costly procedure, an orthogonality regularizer is added to the cost of the loss function. This orthogonality regularizer is given by

$$L_o(\mathbf{U}, \mathbf{V}) = \frac{1}{r^2} \left( \|\mathbf{U}^\top \mathbf{U} - \mathbf{I}\|_F^2 + \|\mathbf{V}^\top \mathbf{V} - \mathbf{I}\|_F^2 \right). \quad (9.1)$$

This assures that the matrices will be approximately semi-unitary. Note that they can not be orthogonal since they are not square.

As part of the sparsification, another regularizer is used for sparsifying the singular values. This regularizer was tested using both the  $L^1$  norm and a *Hoyer regularizer*, given by

$$L^H(\mathbf{s}) = \frac{\|\mathbf{s}\|_1}{\|\mathbf{s}\|_2}. \quad (9.2)$$

The results given by Figure 9 show that while the  $L^1$  regularizer doesn't match the performance/accuracy tradeoff of the hoyer regularizer, allowing for a higher accuracy loss can result in better performance by the  $L^1$  regularizer. I.e. extremely large compression rates which result in higher accuracy loss can provide a more generous reduction in #FLOPs under  $L^1$ , whereas having a more moderate compression and a lower accuracy loss has better reduction of #FLOPs under  $L^H$ .

The overall objective function is finally given by

$$L(\mathbf{U}, \mathbf{s}, \mathbf{V}) = L_T(\text{diag}(\sqrt{\mathbf{s}}) \mathbf{V}^\top, \mathbf{U} \text{diag}(\sqrt{\mathbf{s}})) + \lambda_o \sum_{l=1}^D L_o(\mathbf{U}^{(l)}, \mathbf{V}^{(l)}) + \lambda_s \sum_{l=1}^D L_s(\mathbf{s}^{(l)}), \quad (9.3)$$

where  $L_T$  is the training loss computed with decomposed layers,  $L_o$  is the orthogonality loss from Equation 9.1,  $L_s$  is the sparsity-inducing regularization loss (either  $L^1$  or  $L^H$ ), and  $\lambda_o, \lambda_s$  are hyperparameters for controlling the effect of the objective functions on the overall loss. The paper demonstrates that the change in accuracy from changing  $\lambda_o = 1.0$  to  $\lambda_o = 0.0$  is about  $-2.0\%$ , showing that the orthogonality causes a definite change in performance.

As for pruning, an energy threshold  $e \in [0, 1]$  is a hyperparameter which for each layer, a set  $\mathbb{K}$  is determined containing the largest number of singular values s.t.

$$\sum_{j \in \mathbb{K}} s_j^2 \leq e \sum_{i=1}^r s_i^2. \quad (9.4)$$

The values in  $\mathbb{K}$  are then pruned by setting each  $s_j = 0$ .

**Fill in some details on channel-wise and spatial-wise decomposition because I have no idea what either of these things are.**

## Chapter 10

# Stabilizing Gradients for Deep Neural Networks via Efficient SVD Parameterization [14]

**Abstract.** *Vanishing and exploding gradients are two of the main obstacles in training deep neural networks, especially in capturing long range dependencies in recurrent neural networks (RNNs). In this paper, we present an efficient parametrization of the transition matrix of an RNN that allows us to stabilize the gradients that arise in its training. Specifically, we parameterize the transition matrix by its singular value decomposition (SVD), which allows us to explicitly track and control its singular values. We attain efficiency by using tools that are common in numerical linear algebra, namely Householder reflectors for representing the orthogonal matrices that arise in the SVD. By explicitly controlling the singular values, our proposed Spectral-RNN method allows us to easily solve the exploding gradient problem and we observe that it empirically solves the vanishing gradient issue to a large extent. We note that the SVD parameterization can be used for any rectangular weight matrix, hence it can be easily extended to any deep neural network, such as a multi-layer perceptron. Theoretically, we demonstrate that our parameterization does not lose any expressive power, and show how it potentially makes the optimization process easier. Our extensive experimental results also demonstrate that the proposed framework converges faster, and has good generalization, especially in capturing long range dependencies, as shown on the synthetic addition and copy tasks, as well as on MNIST and Penn Tree Bank data sets.*

The primary contributions of the paper is in

- parameterizing transition matrices through SVD and householder reflectors;
- applying SVD parameterization to RNNs to exert constraints on their transition matrices, and

- verifying their modification outperforms the original RNN networks.

They also theoretically validate their methods as well as show that it alleviates the gradient vanishing/exploding problem.

## 10.1 SVD and Orthogonality with Householder Reflectors

They define the SVD of a transition matrix  $\mathbf{W} \in \mathbb{R}^{n \times n}$  as given by  $\mathbf{W} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top$  where  $\boldsymbol{\Sigma}$  is the diagonal matrix of singular values and  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times n}$  are orthogonal.

**Definition 10.1.** Given a vector  $\mathbf{u} \in \mathbb{R}^k$ ,  $k \leq n$ , the Householder Reflector  $\mathcal{H}_k^n(\mathbf{u})$  is

$$\mathcal{H}_k^n(\mathbf{u}) = \begin{cases} \begin{bmatrix} \mathbf{I}_{n-k} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_k - 2 \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2} \end{bmatrix}, & \mathbf{u} \neq \mathbf{0} \\ \mathbf{I}_n, & \text{otherwise.} \end{cases} \quad (10.1)$$

When  $n$  is well-defined, the previous function is shorthanded to  $\mathcal{H}_k(\mathbf{u})$ .

Let  $\mathbb{O}(n)$  be the set of  $n \times n$  orthogonal matrices. The following maps are established as precursors to the primary theorems of the paper:

$$\begin{aligned} \mathcal{M}_k : \mathbb{R}^k \times \cdots \times \mathbb{R}^n &\rightarrow \mathbb{R}^{n \times n} \\ \mathcal{M}_k(\mathbf{u}_k, \dots, \mathbf{u}_n) &\mapsto \mathcal{H}_n(\mathbf{u}_n) \cdots \mathcal{H}_k(\mathbf{u}_k), \end{aligned} \quad (10.2)$$

and

$$\begin{aligned} \mathcal{M}_{k_1, k_2} : \mathbb{R}^{k_1} \times \cdots \times \mathbb{R}^n \times \mathbb{R}^{k_2} \times \cdots \times \mathbb{R}^n \times \mathbb{R}^n &\rightarrow \mathbb{R}^{n \times n} \\ \mathcal{M}_{k_1, k_2}(\mathbf{u}_{k_1}, \dots, \mathbf{u}_n, \mathbf{v}_{k_1}, \dots, \mathbf{v}_n, \boldsymbol{\sigma}) &\mapsto \\ \mathcal{H}_n(\mathbf{u}_n) \cdots \mathcal{H}_{k_1}(\mathbf{u}_{k_1}) \text{diag}(\boldsymbol{\sigma}) \mathcal{H}_{k_2}(\mathbf{v}_{k_2}) \cdots \mathcal{H}_n(\mathbf{v}_n). \end{aligned} \quad (10.3)$$

**Theorem 10.1.**  $\text{im}(\mathcal{M}_1)$  is the set of all  $n \times n$  orthogonal matrices.

**Theorem 10.2.**  $\text{im}(\mathcal{M}_{1,1})$  is the set of all  $n \times n$  real matrices.

**Theorem 10.3.**  $\text{im}(\mathcal{M}_{k_1, k_2})$  includes the set of all orthogonal  $n \times n$  matrices if  $k_1 + k_2 \leq n + 2$ .

These definitions can be extended to the case of  $\mathbf{W} \in \mathbb{R}^{m \times n}$  by taking

$$\mathbf{W} = \mathbf{U}(\boldsymbol{\Sigma}|\mathbf{0})(\mathbf{V}_L|\mathbf{V}_R)^\top = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}_L^\top, \quad (10.4)$$

where  $\mathbf{U} \in \mathbb{R}^{m \times m}$ ,  $\boldsymbol{\Sigma} \in \text{diag}(\mathbb{R}^m)$ , and  $\mathbf{V}_L \in \mathbb{R}^{n \times m}$ . Equation 10.3 can then be extended as

$$\begin{aligned} \mathcal{M}_{k_1, k_2}^{m,n} : \mathbb{R}^{k_1} \times \cdots \times \mathbb{R}^m \times \mathbb{R}^{k_2} \times \cdots \times \mathbb{R}^n \times \mathbb{R}^{\min(m,n)} &\rightarrow \mathbb{R}^{m \times n} \\ \mathcal{M}_{k_1, k_2}^{m,n}(\mathbf{u}_{k_1}, \dots, \mathbf{u}_m, \mathbf{v}_{k_1}, \dots, \mathbf{v}_n, \boldsymbol{\sigma}) &\mapsto \\ \mathcal{H}_m^m(\mathbf{u}_m) \cdots \mathcal{H}_{k_1}^m(\mathbf{u}_{k_1}) \hat{\boldsymbol{\Sigma}} \mathcal{H}_{k_2}^n(\mathbf{v}_{k_2}) \cdots \mathcal{H}_n^n(\mathbf{v}_n). \end{aligned} \quad (10.5)$$

where  $\hat{\boldsymbol{\Sigma}} = (\text{diag}(\boldsymbol{\sigma})|\mathbf{0})$  if  $m < n$  and  $(\text{diag}(\boldsymbol{\sigma})|\mathbf{0})^\top$  otherwise. This leads to the following lemma and consequently theorem:

**Lemma 10.1.** Given  $\{v_i\}_{i=1}^n$ , define  $V^{(k)} = \mathcal{H}_n^n(\mathbf{v}_n) \cdots \mathcal{H}_k^n(\mathbf{v}_k)$  for  $k \in [n]$ . We have:

$$V_{*,i}^{(k_1)} = V_{*,i}^{(k_2)}, \forall k_1, k_2 \in [n], i \leq \min(n - k_1, n - k_2). \quad (10.6)$$

**Theorem 10.4.** If  $m \leq n$ ,  $\text{im}(\mathcal{M}_{1,n-m+1}^{m,n})$  is the set of all  $m \times n$  matrices; else the image of  $\text{im}(\mathcal{M}_{m-n+1,1}^{m,n})$  is the set of all  $m \times n$  matrices.

Write up notes on the analysis of this paper.

## Chapter 11

# Federated Learning: Strategies for Improving Communication Efficiency [6]

This paper outlines a distributed learning method which operates on the basis that most the devices involved in the learning process are heterogenous and have unreliable internet connectivity. The synchronized federated learning algorithm generally follows these steps:

1. a subset of existing clients is selected, each of which downloads the current model;
2. each client in the subset computes an updated model based on their local data;
3. the model updates are sent from the selected clients to the server; and finally
4. the server aggregates these models (typically by averaging) to construct an improved global model.

The main contribution deals with deriving a method of federated learning which doesn't require full model updates being sent back and forth between the client and server, for which they derive two primary approaches:

**Structured updates** Updates are directly learned from a restricted space that can be parameterized using a smaller number of variables.

**Sketched updates** A full model is learned and then compressed prior to sending it to the server.

The general scheme is to learn a real matrix  $\mathbf{W} \in \mathbb{R}^{d_1 \times d_2}$  from data stored across a large number of clients. This is done by sending a current model  $\mathbf{W}_t$  to a subset

$S_t$  of  $n_t$  clients. Each client updates their local model  $\mathbf{W}_t^i$  (usually via SGD or some learning algorithm) so that its update is given by  $\mathbf{H}_t^i \triangleq \mathbf{W}_t^i - \mathbf{W}_t$ . The next global update is then given by

$$\mathbf{W}_{t+1} = \mathbf{W}_t + \eta_t \mathbf{H}_t, \quad \mathbf{H}_t = \frac{1}{n_t} \sum_{i \in S_t} \mathbf{H}_t^i, \quad (11.1)$$

where  $\eta_t$  is the learning rate.

## 11.1 Structured Updates

*Structured updates* send back and forth full-sized models based on a predefined structure. By doing so, there is little loss of information over the communication channel but potential for higher data throughput.

### 11.1.1 Low Rank

We enforce each  $\mathbf{H}_t^i \in \mathbb{R}^{d_1 \times d_2}$  to be a low rank matrix of rank at most  $k$  for a fixed  $k$  by letting  $\mathbf{H}_t^i = \mathbf{A}_t^i \mathbf{B}_t^i$  where  $\mathbf{A}_t^i \in \mathbb{R}^{d_1 \times k}$  and  $\mathbf{B}_t^i \in \mathbb{R}^{k \times d_2}$ . This matrix is further compressed by representing  $\mathbf{A}_t^i$  as a random matrix defined by a seed and only optimizing  $\mathbf{B}_t^i$ , saving a factor of  $d_1/k$  parameters in communication.  $\mathbf{A}_t^i$  is generated afresh each round and for each client independently.

### 11.1.2 Random Mask

We restrict each  $\mathbf{H}_t^i$  to be sparse based on a predefined sparsity pattern with respect to a random seed (similar to low rank). The sparsity pattern is generated afresh each round and for each client independently.

## 11.2 Sketched Updates

*Sketched updates* place higher emphasis on communication cost rather than fully structured models.

### 11.2.1 Subsampling

We only communicate a matrix  $\hat{\mathbf{H}}_t^i = \mathbf{H}_t^i \odot \mathbf{M}$  where  $\mathbf{M}$  is a sparse mask and is generated afresh each round and for each client independently. This can be done s.t.  $\mathbb{E} [\hat{\mathbf{H}}_t] = \mathbf{H}_t$ .

### 11.2.2 Probabilistic Quantization

Let  $\mathbf{h} = (h_1, \dots, h_{d_1 \times d_2}) = \text{vec}(\mathbf{H}_t^i)$ ,  $h_{\max} = \max_j h_j$ , and  $h_{\min} = \min_j h_j$ . We can apply  $b$ -bit quantization by dividing  $[h_{\min}, h_{\max}]$  into  $2^b$  intervals. Then if  $h_i$

falls in the interval  $[h_{\min}^i, h_{\max}^i]$ , then the compressed update for index  $i$  is defined by

$$\tilde{h}_i = \begin{cases} h_{\max}^i, & \text{with probability } \frac{h_i - h_{\min}^i}{h_{\max}^i - h_{\min}^i} \\ h_{\min}^i, & \text{with probability } \frac{h_{\max}^i - h_i}{h_{\max}^i - h_{\min}^i}. \end{cases} \quad (11.2)$$

Using this, we have that  $\mathbb{E} [\tilde{\mathbf{h}}] = \mathbf{h}$ .

Error bounds on this method can be improved by a factor of  $\mathcal{O}(d/\log d)$  by first performing a random rotation on  $\mathbf{h}$ , quantizing, sending the information to the server, then performing the inverse rotation. The paper states that a structured rotation matrix, the product of a Walsh-Hadamard and binary diagonal entry matrix, can have a computational complexity of  $\mathcal{O}(d)$  for generating the matrix and  $\mathcal{O}(d \log d)$  for applying the matrix.

## Chapter 12

# Scalable Adaptive Stochastic Optimization Using Random Projections [7]

This paper uses a similar notation to that of Chapter 7 in that we consider the function

$$F_t = f_t(\beta) + \phi(\beta) \quad (12.1)$$

where  $f_t$  is a convex loss function,  $\phi$  is a regularizer, and  $\beta \in \mathbb{R}^p$  is the set of parameters.

**Definition 12.1.** The *regret* after  $T$  rounds is defined as

$$R(T) = \sum_{t=1}^T F_t(\beta_t) - \sum_{t=1}^T F_t(\beta^{\text{opt}}), \quad (12.2)$$

where  $\beta_t$  is the set of parameters at round  $t$  and  $\beta^{\text{opt}}$  is a fixed optimal predictor.

The paper proceeds to define  $\mathbf{g}_t \in \nabla f_t(\beta_t)$  as a *subgradient* of the loss function, which we may view as a gradient obtained from a minibatch of our dataset. Standard, first-order methods use a fixed learning rate  $\eta$  to move in the opposite direction of  $\mathbf{g}_t$ , but their method builds on second-order methods which use an adaptive learning rate which is controlled by a time-varying *proximal term*, i.e. a function determining the adaptive learning rate. Given  $\mathbf{G}_t = \sum_{i=1}^t \mathbf{g}_i \mathbf{g}_i^\top$  and  $\mathbf{H}_t = \delta \mathbf{I}_p + \mathbf{G}_t^{1/2}$  The proximal terms for ADAFULL and ADAGRAD are given by

$$\psi_t(\beta) = \frac{1}{2} \langle \beta, \mathbf{H}_t \beta \rangle \quad \text{and} \quad \psi_t(\beta) = \frac{1}{2} \left\langle \beta, \left( \delta \mathbf{I}_p + \text{diag}(\mathbf{G}_t)^{1/2} \right) \beta \right\rangle \quad (12.3)$$

respectively. When  $p$  is large, ADAFULL may be computationally infeasible, yet ADAGRAD is unable to capture dependencies between coordinates in gradient terms.

**Definition 12.2.** *Effective rank* is defined for a matrix  $\Sigma$  as

$$r(\Sigma) = \frac{\text{tr}(\Sigma)}{\|\Sigma\|}. \quad (12.4)$$

$\Sigma$  is considered to have *low effective rank* when  $r(\Sigma) \leq \text{rank}(\Sigma) \leq p$ .

*Remark.* Low effective rank is a common trademark of high-dimensional data, as there is generally a lot of sparsity.

**Definition 12.3.** *Random projections* are low-dimensional embeddings  $\Pi : \mathbb{R}^p \rightarrow \mathbb{R}^\tau$  which preserve - up to a small distortion - the geometry of a subspace of vectors.

**Definition 12.4.** The *Subsampled Randomized Fourier Transform* is a *structured* random projection given by

$$\Pi = \sqrt{p/\tau} S \Theta D \quad (12.5)$$

where

1.  $S \in \mathbb{R}^{\tau \times p}$  is a subsampling matrix,
2.  $D \in \mathbb{R}^{p \times p}$  is a diagonal matrix whose entries are drawn independently from  $\{-1, 1\}$ ,
3.  $\Theta \in \mathbb{R}^{p \times p}$  is a unitary discrete fourier transform matrix.

They propose the following algorithms.

---

**Algorithm 2** ADA-LR

---

**Require:**  $\eta > 0, \delta \geq 0, \tau$

**for**  $t = 1 \dots T$  **do**

    Receive  $\mathbf{g}_t = \nabla f_t(\beta_t)$

$\mathbf{G}_t = \mathbf{G}_{t-1} + \mathbf{g}_t \mathbf{g}_t^\top$  ▷ Construct additive gram matrix

    Project  $\tilde{\mathbf{G}}_t = \mathbf{G}_t \Pi$

$\mathbf{Q}\mathbf{R} = \tilde{\mathbf{G}}_t$  ▷ QR Decomposition to get the basis vectors  $\mathbf{Q}$

$\mathbf{B} = \mathbf{Q}^\top \mathbf{G}_t$  ▷ Project onto new basis

$\mathbf{U}, \Sigma, \mathbf{V} = \mathbf{B}$  ▷ SVD

$\beta_{t+1} = \beta_t - \eta \mathbf{V} (\Sigma^{1/2} + \delta \mathbf{I})^{-1} \mathbf{V}^\top \mathbf{g}_t$

**end for**

---

“

.....

## Part V

# Finance

# Chapter 13

## An Introduction to Quantitative Finance [2]

It is important to note that I only present the conclusions to propositions presented in this book rather than deriving the proofs themselves. This is especially the case whenever I list out the value of any derivative: these values are derived in the book using various proofs and are *not* defined to have these values in the first place. For the most part, any derivative with maturity  $T$  has a definitive value at maturity, but is a random variable up until that point. The book spends a good amount of time providing deterministic values for these r.v.s, and all that is generally listed here are these deterministic values.

### 13.1 Interest Rates

**Definition 13.1.** A *notional* or *principal* is an initial deposit or value  $N$ .

**Definition 13.2.** An *interest rate*  $r$  is the rate at which a value is increased according to a specified frequency of time.

Suppose we have an account with interest rate  $r$  and notional  $N$ . The value of the account at time  $T$  is

$$N_T = Ne^{rT} \tag{13.1}$$

Throughout finance, there exist discrete and continuous analogs for most equations due the discrete case being more realistic while the continuous case is far more mathematically efficient. That being said, discrete analogs will be used sparsely and we will primarily depend on continuous equations. The discrete versions can be derived from the continuous versions by noting that

$$e^{rT} = \left(1 + \frac{r_m}{m}\right)^{mT} \tag{13.2}$$

where  $T$  is commonly denoted as the time in years,  $m$  is the frequency the interest is compounded, and  $r_m$  is the equivalent interest rate with the discrete frequency.

Table 13.1: Daycount conventions and accrual factors

Daycount convention	$\alpha$ for 16 December 2011 - 16 March 2012
act/365	91/365
act/act	15/365 + 76/366
act/360	91/360
30/365	1/4

$r_m$  can then be recovered from  $r$  using

$$r_m = m \left( e^{\left(\frac{r}{m}\right)} - 1 \right). \quad (13.3)$$

**Definition 13.3.** A *money market account* is an account compounded continuously at rate  $r$  with notional 1. I.e.  $M_0 = 1$  and  $M_t = e^{rt}$ .

**Definition 13.4.** A *zero coupon bond* (ZCB) with maturity  $T$  is an asset that pays 1 at time  $T$  (and nothing else). The value of a ZCB at time  $t$  for a continuously compounded rate  $r$  is denoted as

$$Z(t, T) = e^{-r(T-t)} \quad (13.4)$$

where  $Z(T, T) = 1$  by definition. The values  $Z(t, T)$  with  $0 \leq t \leq T$  are known as *discount factors* or *present values*.

The intuition behind the present values for a ZCB derives from the fact that if we have two portfolios, one a ZCB with maturity  $T$  and another that contains cash which will accumulate interest to be worth 1 at time  $T$ , then both will be worth the same value at each  $t$  leading up to  $T$  and thereafter.

**Definition 13.5.** An *annuity* a series of fixed cashflows at specified times  $T_i = 1, \dots, n$ . The value at time  $t$  is given by

$$V = C \sum_{i=1}^n Z(t, T_i). \quad (13.5)$$

When we substitute in Equation 13.4, we notice that this is just a geometric series. Because of this, if an annuity is payed once a year for  $M$  years with rate  $r$ , this equates to

$$V = \sum_{i=1}^n \frac{1}{(1+r)^i} = \frac{1}{r} \left( 1 - \frac{1}{(1+r)^M} \right) \quad (13.6)$$

**Definition 13.6.** The *accrual factor* is (under the discrete analogy) the frequency which interest is compounded, and is denoted as  $\alpha$ .

In Equations 13.2 and 13.3, this is the same as taking  $\alpha = 1/m$ . Due to the year being uneasily divided by simple values of  $\alpha$ , several *daycount conventions* have been made for different settings. We list a sample of these in Table 13.1. The daycount conventions are seemingly deceptive in that they hide various details

behind terminology. For example, two cases are  $30/x$  and  $x/360$ , which respectively denote that each month in the period should be counted as only having 30 days and each month in a year should be counted as only having 30 days. I want to emphasize the weird convention because when we encounter a fraction in the work such as  $30/\text{act}$ , this means to say that if we specify a period of three months, then this fraction actually denotes

$$\frac{30 * 3}{m_1 + m_2 + m_3} \quad (13.7)$$

where each  $m_i$  is the number of actual days in the specified month. Like wise, if we have a daycount  $\text{act}/\text{act}$ , this denotes the number specifies the number of actual days in the period divided by the actual days in the year. If we have a leap year, then the denominator contains 366, whereas if February is included in the daycount, then we may have 28 or 29 days dependent on the year. I hope this helps to explain it better, if not then *oof*.

A few specific cases arise but are only defined in the exercises. I will list those here.

**Definition 13.7.** *Semi-bond* is the US dollar interest rate which is semi-annual compounding with  $30/360$  daycount and is denoted as  $\gamma_{SB}$ .

**Definition 13.8.** *Annual money* is the US dollar interest rate which is annual compounding with  $\text{act}/360$  daycount and is denoted as  $\gamma_{AM}$ .

If we know the rate  $r_\alpha$  and accrual factor  $\alpha$  of a quoted interest rate over a specified period of time, and we are interested in the equivalent rate  $r_\beta$  under a different accrual factor  $\beta$  over the same period of time, then this information can be retrieved by making note of the equivalency

$$r_\alpha \alpha = r_\beta \beta. \quad (13.8)$$

This makes conversions simple as demonstrated by the following example.

**Example 13.1.** If on 16 December 2011 the  $\text{act}/365$  rate for three months is 5%, the  $\text{act}/360$  rate is

$$r_\beta = r_\alpha \frac{\alpha}{\beta} = 5\% \frac{\text{act}/365}{\text{act}/360} = 5\% \frac{360}{365} = 4.9315\%. \quad (13.9)$$

**Example 13.2.** The conversion between semi-bonds and annual money is given by

$$\gamma_{SB} \frac{30}{360} = \gamma_{AM} \frac{\text{act}}{360}, \quad (13.10)$$

so that

$$\gamma_{SB} = \gamma_{AM} \frac{\text{act}}{30}, \quad (13.11)$$

and likewise

$$\gamma_{AM} = \gamma_{SB} \frac{30}{\text{act}}. \quad (13.12)$$

**Definition 13.9.** A *stock* or *share* is an asset giving ownership in a fraction of a company. The price at time  $t$  of a stock is denoted by  $S_t$ . The current known price is called the *spot*.

**Definition 13.10.** A *fixed rate bond* with coupon  $c$  and notional  $N$  is an asset that pays a coupon  $cN$  at a specified frequency (usually a year) and  $N$  at its maturity date  $T$ . A *floating rate bond* is defined similarly with a variable interest rate.

## 13.2 Forward Contracts and Forward Prices

**Definition 13.11.** A *derivative contract (derivative)* is a financial contract between two parties whose value derives from the value of another variable.

*Remark.* Derivatives are often defined by the payout at a specified time. E.g. if  $S_T$  indicates the total snowfall in inches at time  $T$  then the weather derivative  $g(S_T)$  can be defined by

$$g(S_T) = I\{S_T > 50\} = \begin{cases} \$1 & \text{if } S_T > 50 \\ 0 & \text{otherwise,} \end{cases} \quad (13.13)$$

where both  $S_T$  and  $g(S_T)$  are random variables with unknown values until  $T = 1$ .

**Definition 13.12.** A *forward contract (forward)* is an agreement between two counterparties to trade a specific asset at maturity  $T$  with delivery price  $K$ . At a time  $t \leq T$ , the buyer is *long* the contract while the seller is *short* the contract.

The value at time  $t \leq T$  of being long a forward contract is given by  $V_K(t, T)$ , where by construction we have that  $V_K(T, T) = S_T - K$ .

**Definition 13.13.** The *forward price*  $F(t, T)$  at current time  $t \leq T$  is the delivery price  $K$  s.t.  $V_K(t, T) = 0$ .

*Remark.* The forward price can be seen as denoting the spot price at maturity, so in essence you will be agreeing at time  $t \leq T$  to trade the asset at time  $T$  for its actual price. If you had  $K < F(t, T)$  then the person long the contract has a better deal since they would pay less than the appreciated value (i.e. you could just sell the asset for profit), whereas  $K > F(t, T)$  provides the better deal for the short since they would receive more than the appreciated value.

Table 13.2: Forward prices for various assets

Underlying	Forward price
Asset paying no income	$S_t e^{r(T-t)}$
Asset paying known income $I$	$(S_t - I) e^{r(T-t)}$
Asset paying dividends at rate $q$	$S_t e^{(r-q)(T-t)}$
Foreign exchange	$X_t e^{(r_s - r_f)(T-t)}$

The forward price is decided under various conditions, those being listed in Table 13.2. Given the forward price, we derive the value of being long a forward contract as

$$V_K(t, T) = (F(t, T) - K)e^{-r(T-t)}. \quad (13.14)$$

As a final note, this chapter begins to discuss proof methods for deriving the values of derivatives. It's fairly instructional to review these as they form the basis for evaluating many assets, so we shall include the definitions and some proof examples below.

**Definition 13.14.** A *replication proof* is one where we prove the value of an asset by demonstrating that if two portfolios always have the same value at time  $T$  and if neither add/subtract anything of non-zero value between  $t$  and  $T$ , then they must have the same value at  $t \leq T$ .

**Definition 13.15.** *Arbitrage* is a situation where, starting with an empty portfolio, a sequence of simple market transactions will end up with a positive portfolio value at time  $T$ .

**Definition 13.16.** A *no-arbitrage proof* uses the assumption that no arbitrage situations exist, i.e. all market transactions have an equivalent exchange.

*Remark.* Most no-arbitrage proofs are constructed similarly to a proof by contradiction.

It's often informative constructing no-arbitrage proofs since they force you to think about situations in which you can exploit an arbitrage if detected. We demonstrate no-arbitrage proofs and replication proofs in the following examples.

**Example 13.3.** Let  $S_t$  be the current stock of a price paying no income. Let  $r_{BID}$  be the interest rate at which one can invest/lend money, and  $r_{OFF}$  be the interest rate at which one can borrow money,  $r_{BID} \leq r_{OFF}$ . Both rates are continuously compounded. Using arbitrage arguments, find upper and lower bounds for the forward price of the stock for a forward contract with maturity  $T > t$ .

*Solution.* Assume that  $F(t, T) > S_t e^{r_{OFF}(T-t)}$ . At time  $t$ , we could short a forward contract with forward price  $F(t, T)$  (which is done at no upfront cost since  $V_{F(t, T)}(t, T) = 0$ ), take out a loan for  $S_t$  at rate  $r_{OFF}$  for time  $T-t$ , and purchase the asset for  $S_t$ . Then at time  $T$ , we would execute the forward contract, selling the stock for  $F(t, T)$  and paying back the loan with appreciated value  $S_t e^{r_{OFF}(T-t)}$ , leaving us with  $F(t, T) - S_t e^{r_{OFF}(T-t)} > 0$ . This provides an upper bound under the assumption of no-arbitrage.

Assume that  $F(t, T) < S_t e^{r_{BID}(T-t)}$ . At time  $t$ , we could long a forward contract with forward price  $F(t, T)$ , borrow the asset  $S_t$ , sell the asset and invest/lend the remains at the rate  $r_{BID}$  for time  $T-t$ . Then at time  $T$ , we receive  $S_t e^{r_{BID}(T-t)}$  at the maturity of the deposit, which we can use  $F(t, T)$  of to execute the forward contract and return the underlying asset. This leaves us with  $S_t e^{r_{BID}(T-t)} - F(t, T) > 0$ . This provides a lower bound under the assumption of no-arbitrage.

Our bounds then come out to

$$S_t e^{r_{BID}(T-t)} \leq F(t, T) \leq S_t e^{r_{OFF}(T-t)}. \quad (13.15)$$

## Chapter 14

# A Random Walk Down Wall Street [8]

# **Part VI**

# **Differential Privacy**

# Bibliography

- [1] Gagan Aggarwal et al. “Approximation Algorithms for k-Anonymity”. In: *Proceedings of the International Conference on Database Theory (ICDT 2005)*. Nov. 2005. URL: <http://ilpubs.stanford.edu:8090/645/>.
- [2] S. Blyth. *An Introduction to Quantitative Finance*. OUP Oxford, 2013. ISBN: 9780199666591. URL: <https://books.google.com/books?id=SXbcAAAAQBAJ>.
- [3] Michel X. Goemans and David P. Williamson. “Improved Approximation Algorithms for Maximum Cut and Satisfiability Problems Using Semidefinite Programming”. In: *J. ACM* 42.6 (Nov. 1995), pp. 1115–1145. ISSN: 0004-5411. DOI: 10.1145/227683.227684. URL: <https://doi.org/10.1145/227683.227684>.
- [4] Ian Goodfellow. *Efficient Per-Example Gradient Computations*. 2015.
- [5] Pinar Heggernes. *Treewidth, partial k-trees, and chordal graphs*. Sept. 2006.
- [6] Jakub Konečný et al. “Federated Learning: Strategies for Improving Communication Efficiency”. In: *CoRR* abs/1610.05492 (2016). arXiv: 1610.05492. URL: <http://arxiv.org/abs/1610.05492>.
- [7] Gabriel Krummenacher et al. *Scalable Adaptive Stochastic Optimization Using Random Projections*. 2016. arXiv: 1611.06652 [stat.ML].
- [8] Burton. G. Malkiel. *A Random Walk Down Wall Street*. Norton, New York, 1973.
- [9] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. 2e. New York, NY, USA: Springer, 2006.
- [10] Jiang Qian et al. “Understanding gradient clipping in incremental gradient methods”. In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2021, pp. 1504–1512.
- [11] Hidenori Tanaka et al. “Pruning neural networks without any data by iteratively conserving synaptic flow”. In: *CoRR* abs/2006.05467 (2020). arXiv: 2006.05467. URL: <https://arxiv.org/abs/2006.05467>.
- [12] Roberto Tempo and Hideaki Ishii. “Monte Carlo and Las Vegas Randomized Algorithms for Systems and Control\*: An Introduction”. In: *European Journal of Control* 13.2 (2007), pp. 189–203. ISSN: 0947-3580. DOI: <https://doi.org/10.3166/ejc.13.189-203>. URL: <https://www.sciencedirect.com/science/article/pii/S0947358007708191>.

- [13] Huanrui Yang et al. “Learning Low-rank Deep Neural Networks via Singular Vector Orthogonality Regularization and Singular Value Sparsification”. In: *CoRR* abs/2004.09031 (2020). arXiv: 2004 . 09031. URL: <https://arxiv.org/abs/2004.09031>.
- [14] Jiong Zhang, Qi Lei, and Inderjit Dhillon. “Stabilizing Gradients for Deep Neural Networks via Efficient SVD Parameterization”. In: *Proceedings of the 35th International Conference on Machine Learning*. Ed. by Jennifer Dy and Andreas Krause. Vol. 80. Proceedings of Machine Learning Research. PMLR, July 2018, pp. 5806–5814. URL: <https://proceedings.mlr.press/v80/zhang18g.html>.
- [15] Xiangyu Zhang et al. “Accelerating Very Deep Convolutional Networks for Classification and Detection”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 38.10 (2016), pp. 1943–1955. DOI: 10.1109/TPAMI.2015.2502579.