

## 論文輪講

# Federated Optimization: Distributed Machine Learning for On-Device Intelligence arXiv:1610.02527, Jakub Konečný, et al., 2016

杉浦 圭祐

慶應義塾大学理工学部情報工学科 松谷研究室

May 16, 2019

# 目次

- 1 問題設定
- 2 基本的な最適化アルゴリズム
- 3 ランダム化された最適化アルゴリズム
- 4 Federated Learning のための最適化アルゴリズム
- 5 最適化アルゴリズムの比較
- 6 結論

## 1 問題設定

- 問題の定式化

- 解くべき問題は次のように定式化される

$$\min_{\mathbf{w} \in \mathbb{R}^D} f(\mathbf{w}), \quad f(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{w}) \quad (1)$$

- 入出力のデータを  $\{\mathbf{x}_i, y_i\}_{i=1}^N$ 、損失関数を  $f_i(\mathbf{w})$  とする
- 具体的には、以下のような問題が考えられる

線形回帰：

$$f_i(\mathbf{w}) = \frac{1}{2} (\mathbf{x}_i^T \mathbf{w} - y_i)^2, \quad y_i \in \mathbb{R} \quad (2)$$

ロジスティック回帰：

$$f_i(\mathbf{w}) = -\log(1 + \exp(-y_i \mathbf{x}_i^T \mathbf{w})), \quad y_i \in \{-1, 1\} \quad (3)$$

サポートベクタマシン：

$$f_i(\mathbf{w}) = \max\{0, 1 - y_i \mathbf{x}_i^T \mathbf{w}\}, \quad y_i \in \{-1, 1\} \quad (4)$$

- 問題の定式化

- 上記は凸最適化問題である
- 線形回帰、ロジスティック回帰、サポートベクタマシンより複雑なモデルにも適用可能
  - ⇒ 条件付き確率場や、ニューラルネットワーク
  - ⇒ 損失関数  $f_i(\mathbf{w})$  が非凸で、複雑な形状をしている場合
- データ数  $N$  が大き過ぎて、単一のノードで学習を進めるのは不可能
  - ⇒ 分散処理が必要 (データが複数の箇所に分散し、複数の相互接続されたノードで学習を行う)
  - ⇒ ノード間での通信時間がボトルネックとなる可能性
  - ⇒ 並列計算の利点を活かすために、モデルの学習を、単一ノードでの計算に適した簡単な問題に分割する必要がある

- 最新 (State-of-the-art) の最適化手法
  - シーケンシャルであるため、並列処理には向かない
  - 各イテレーションでの処理は非常に高速だが、イテレーションの実行を何度も繰り返す必要がある
    - ⇒ 各イテレーションの実行後に、複数のノード間で通信を行うと、性能が極端に落ちる
    - ⇒ 各イテレーションの実行時間より、ノード間の通信時間の方が遥かに大きい

- 従来手法における仮定

- 1 データは各ノードに均等に分散している

- 2 ノード数  $K$  に比べて、1つのノードが保持しているデータ数の平均  $N/K$  の方が非常に大きい ( $K \ll N/K$ )  
⇒ 大規模なデータセンタに、データが格納されている想定

- 3 各ノードが、分布をよく表現するデータを保持している  
⇒ 各ノードが、**独立同分布** (IID) 標本を持っているという前提

- ⇒ 実際には、ノードの地理的な位置によって、各ノードが持つデータが、クラスタに分かれている可能性

- ⇒ 各ノードが持つデータが時間的に変動し、ある時点では他のノードと似たようなデータを保持している可能性

- ⇒ あるノードに頻繁に現れる特徴が、他のノードでは全く現れない可能性

- 従来手法における仮定

- 従来手法では、上記 3 つの仮定が成立する
  - ⇒ Federated Learning では、これらの**仮定を全く置かない**
- 従来手法では、最初に、デバイス上のデータを中央のノード (データセンタ等) に集める
  - ⇒ 集めたデータをランダムにシャッフルし、複数の計算ノードに均等な数だけ配分して、学習させる
- Federated Learning では、デバイス上のデータを中央のノードに送信しない
  - ⇒ 各デバイスと中央のノードとの通信量が大幅に削減
  - ⇒ ユーザのプライバシーを保護し、セキュリティを向上させる
  - ⇒ データがデバイス上にしか存在しないので、中央のノードが攻撃されても、ユーザのデータが漏洩する危険性がない (攻撃されそうな箇所の候補が減る)



- Federated Learning での問題設定

- Federated Learning では、次のような現実的な仮定を置く

- Massively Distributed**: データは、多数のノードに分散  
⇒ ノード数  $K$  と、1つのノードが保持しているデータ数の平均  $N/K$  を比較したとき、 $N/K \ll K$  となる可能性
- Non-IID**: 各ノードが保持するデータは、異なる分布からサンプルされた可能性  
⇒ あるノードが保持しているデータは、データ全体の分布を表現しているとは限らない  
⇒ 各ノードが保持するデータは、独立に同一の分布からサンプルされた (IID) とは、仮定し難い
- Unbalanced**: ノードによって、保持しているデータ数は大きく異なる

- Federated Learning での問題設定

- 上記に加えて、この論文では次のような仮定を置く

- 1 データはスパースである

⇒ ある特徴は、ほんの僅かなノードやデータにしか現れない

- 2 データはモバイル端末上に存在し、プライバシー上の配慮が必要 (Privacy Sensitive)

⇒ 入出力データ  $\{x_i, y_i\}$  はデバイス上で作られる

⇒ 例えば、ユーザが次に入力する単語の予測、ユーザがシェアしそうな写真の予測、ユーザにとってどの通知が重要かの予測

- 3 多数のデバイスが動作するので、事実上無限の計算能力が得られる
  - ⇒ 但し、各デバイスと中央のノード間の、通信コストによって制限される (バンド幅の制限)
  - ⇒ デバイスと中央のノード間の通信を、いかに減らせるかが、性能向上のための鍵となる
  
- 4 差分データ  $\delta \in \mathbb{R}^D$  が、モデルの学習に使用する**唯一の情報**である
  - ⇒ 各デバイスは、1 回の Round につき、モデルのパラメータの差分  $\delta \in \mathbb{R}^D$  を計算 ( $D$  は、モデルのパラメータの次元)
  - ⇒ 差分  $\delta$  が、デバイスから中央のノードにアップロードされる

- Federated Learning で送信される差分データ  $\delta$  について
  - $\delta$  は、ユーザのプライベートな情報を依然として含むかもしれないが、訓練データ  $\{x_i, y_i\}$  に比べれば無視できるほど小さい
  - モデルのパラメータの差分ベクトル  $\delta$  のサイズは、訓練データのサイズには関係ない  
⇒ 訓練データが巨大 (例えばユーザが撮影した動画) であっても、データそのものではなく、差分データのみを中央のノードに送信すればよい  
ため、**通信量が大幅に削減**できる

# 問題設定

- 差分ベクトル  $\delta$  の使途は、元の訓練データ  $\{x_i, y_i\}$  に比べて限られる
  - ⇒ モデルの訓練以外に、殆ど使い道がない
  - ⇒ 各デバイスから送信された差分データは、モデルの訓練に使用した後、破棄してよい
- ⇒ ユーザにとっては、アップロードしたデータが、モデルの訓練にしか使われない (想定外の方法で使われない) ことが分かっているので、安心できる (**プライバシーの保護**につながる)
- ⇒ モデルを訓練する側にとっては、ユーザのデータを保存する際の負担が軽減される
- ⇒ 訓練データであれば、漏洩しないように厳重に管理する必要がある
- ⇔ 差分データであれば、万が一不正にアクセスされても、ユーザの個人情報漏洩することはない

- 論文で提案された訓練アルゴリズムについて
  - Federated Learning のための訓練アルゴリズムを新たに設計した  
⇐ **Massively Distributed**、**Non-IID**、**Unbalanced**
  - 新たな訓練アルゴリズムによって、比較的少ない Round 数 (通信量) で、モデルのパラメータを収束させることができた

## 2 基本的な最適化アルゴリズム

# 基本的な最適化アルゴリズム

- 解くべき問題は、次のように定式化された
  - $D$  はモデルのパラメータの次元数、 $N$  は訓練データ数
  - $\mathbf{w} \in \mathbb{R}^D$  はパラメータベクトル、 $f_i(\mathbf{w})$  は損失関数

$$\min_{\mathbf{w} \in \mathbb{R}^D} f(\mathbf{w}), \quad f(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{w}) \quad (5)$$

- ベースラインアルゴリズム
  - 基本的なアルゴリズムとして、以下を紹介する
    - 勾配降下法 (GD; Gradient Descent)
    - 確率的勾配降下法 (SGD; Stochastic Gradient Descent)



# 基本的な最適化アルゴリズム

- 勾配降下法 (GD; Gradient Descent)
  - パラメータの更新式は次のようになる

$$\mathbf{w}^{t+1} = \mathbf{w}^t - h_t \nabla f(\mathbf{w}^t) \quad (6)$$

- $\nabla f(\mathbf{w}^t)$  は次のように定義される

$$\nabla f(\mathbf{w}^t) \equiv \left. \frac{\partial}{\partial \mathbf{w}} f(\mathbf{w}) \right|_{\mathbf{w}=\mathbf{w}^t} \quad (7)$$

$$= \frac{1}{N} \sum_{i=1}^N \left. \frac{\partial}{\partial \mathbf{w}} f_i(\mathbf{w}) \right|_{\mathbf{w}=\mathbf{w}^t} \quad (8)$$

$$= \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{w}^t) \quad (9)$$

損失関数  $f(\mathbf{w})$  のパラメータ  $\mathbf{w}$  による勾配の、 $\mathbf{w} = \mathbf{w}^t$  における値

- $h_t > 0$  は学習率 (ステップサイズ)

# 基本的な最適化アルゴリズム

- 勾配降下法 (GD; Gradient Descent) の問題点
  - 勾配  $\nabla f(\mathbf{w}^t)$  を求めるためには、 $N$  個の各データについての勾配  $\nabla f_i(\mathbf{w}^t)$  を計算する必要がある
    - ⇒ 1 回のパラメータ更新のために、全データを処理する必要がある
    - ⇒ データ数  $N$  は非常に大きいため、勾配の計算に時間が掛かり過ぎる
  - ⇒ 勾配降下法は、遅すぎて使い物にならない
- モメンタム項を加えることで、アルゴリズムを高速化できる
  - ⇒ 但し、1 回のパラメータ更新のために、全データを処理する必要がある
  - ⇒ モメンタム項を加えても、やはり使い物にならない

# 基本的な最適化アルゴリズム

- 確率的勾配降下法 (SGD; Stochastic Gradient Descent)

- パラメータの更新式は次のようになる

$$\mathbf{w}^{t+1} = \mathbf{w}^t - h_t \nabla f_{i_t}(\mathbf{w}^t) \quad (10)$$

- $\nabla f_{i_t}(\mathbf{w}^t)$  は次のように定義される

$$\nabla f_{i_t}(\mathbf{w}^t) = \left. \frac{\partial}{\partial \mathbf{w}} f_{i_t}(\mathbf{w}) \right|_{\mathbf{w}=\mathbf{w}^t} \quad (11)$$

- $i_t$  は、1 から  $N$  の中から適当に選択されたインデックス  
⇒ 時刻  $t$  では、 $i_t$  番目のデータ  $\{\mathbf{x}_{i_t}, y_{i_t}\}$  とパラメータ  $\mathbf{w}^t$  を用いて、損失関数  $f_{i_t}(\mathbf{w}^t)$  を計算  
⇒ 損失関数の勾配  $\nabla f_{i_t}(\mathbf{w}^t)$  を求めて、パラメータ  $\mathbf{w}$  を勾配の方向に更新

# 基本的な最適化アルゴリズム

- 確率的勾配降下法 (SGD; Stochastic Gradient Descent)

- 1 回のパラメータ更新のためには、1 つのデータ点に対する勾配  $\nabla f_{i_t}(\mathbf{w}^t)$  だけを求めればよい  
⇒ 勾配降下法とは異なり、全データを処理する必要がない

- 1 つのデータ点に対する勾配の期待値は、損失関数  $f(\mathbf{w})$  の勾配の不偏推定量となっている  
⇒  $\mathbb{E}[\nabla f_{i_t}(\mathbf{w})] = \nabla f(\mathbf{w}^t)$  であり、この手法が正当化される

⇒ 実際には、データ点のサンプリングによって、(真の勾配に対して) ノイズが加わった勾配が得られるため、パラメータの収束が遅くなる  
⇒ 学習率 (ステップサイズ)  $h_t$  の設定が重要になる

- 勾配の計算に用いるデータ点  $i_t$  を、毎回ランダムに選ぶのではなく、全データをランダムな順で取り出し、その勾配を求めてパラメータを更新していく手法がある

# 基本的な最適化アルゴリズム

- GD と SGD との比較

- GD ではパラメータの収束が速い  $\Leftrightarrow$  SGD は収束が遅い
- GD では各イテレーションの計算に非常に時間が掛かる
  - $\Leftarrow$  各イテレーションにおいて、全データを処理する必要がある
  - $\Leftrightarrow$  SGD では各イテレーションの計算は高速であり、計算時間はデータ数  $N$  に依存しない
- 今回解こうとしている問題では、パラメータの精度はそこまで高くなくてもよい
  - $\Rightarrow$  SGD で十分である (極端な場合には、全データを 1 回ずつ処理するだけで、パラメータが収束)
  - $\Leftarrow$  GD の場合は、全データを 1 回ずつ処理して、ようやくパラメータを 1 回更新できる

## ③ ランダム化された最適化アルゴリズム

# ランダム化された最適化アルゴリズム

- ランダム化された座標降下法 (RCD; Randomized Coordinate Descent)
  - パラメータの更新式は次のようになる

$$\mathbf{w}^{t+1} = \mathbf{w}^t - h_{j_t} \nabla_{j_t} f(\mathbf{w}^t) \mathbf{e}_{j_t} \quad (12)$$

- $j_t$  は、1 から  $D$  (パラメータの次元数) の中から適当に選択された次元
- $h_{j_t}$  は学習率 (ステップサイズ)、 $\mathbf{e}_{j_t} \in \mathbb{R}^D$  は次元  $j_t$  方向の標準基底ベクトル
- 勾配  $\nabla_{j_t} f(\mathbf{w}^t)$  は次のように定義される

$$\nabla_{j_t} f(\mathbf{w}^t) = \frac{\partial}{\partial w_{j_t}} f(\mathbf{w}^t) \quad (13)$$

$$= \frac{1}{N} \sum_{i=1}^N \left. \frac{\partial}{\partial w_{j_t}} f_i(\mathbf{w}) \right|_{\mathbf{w}=\mathbf{w}^t} \quad (14)$$

$$= \frac{1}{N} \sum_{i=1}^N \nabla_{w_{j_t}} f_i(\mathbf{w}^t) \quad (15)$$

# ランダム化された最適化アルゴリズム

- ランダム化された座標降下法 (RCD; Randomized Coordinate Descent)
  - 最適化問題  $\min_{\mathbf{w} \in \mathbb{R}^D} f(\mathbf{w})$  を、幾つかの部分問題に分割
  - $f(\mathbf{w})$  を、 $j \neq j_t$  であるような変数については固定したまま、ある 1 つの変数  $j_t \in \{1, \dots, D\}$  について最小化する  
⇒ 一度に 1 つの座標を最適化するため、**座標降下法**とよばれる
  - ⇒ 一般に、変数の部分集合について同時に最小化を行うアルゴリズムを、**ブロック座標降下法**という
- 座標降下法は、ある 1 つの変数が、他の変数の最適値に影響を与えない場合に効果を発揮  
⇒ 「深層学習」の 8.7.2 節を参照



# ランダム化された最適化アルゴリズム

- SDCA; Stochastic Dual Coordinate Ascent

- 正則化項  $\psi(\mathbf{w})$  を付加した、以下の最適化問題を考える

$$\min_{\mathbf{w} \in \mathbb{R}^D} f(\mathbf{X}\mathbf{w}) + \psi(\mathbf{w}), \quad f(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N f_i(y_i, \mathbf{x}_i^T \mathbf{w}) \quad (16)$$

- 上記の**双対問題** (Dual problem) は、**Fenchel 双対定理**から、次のようになる ( $K$  はデータ  $\mathbf{x}$  の次元数) [3]

$$- \min_{\mathbf{u} \in \mathbb{R}^N} f^*(\mathbf{u}) + \psi^* \left( -\frac{\mathbf{X}^T \mathbf{u}}{N} \right) \quad (17)$$

- 但し、 $f^*$  と  $\psi^*$  は、それぞれ  $f$  と  $\psi$  の**ルジャンドル変換**である

# ランダム化された最適化アルゴリズム

- SDCA; Stochastic Dual Coordinate Ascent

- $f(x)$  のルジャンドル変換  $f^*(y)$  とは、関数  $f(x)$  の変数を、その微分  $y = x'$  に置き換えた関数であり、次のように定義される (関数  $f(x)$  を、その傾きの情報から捉えた関数)

$$f^*(y) = \sup_x \{x^T y - f(x)\} \quad (18)$$

- 双対問題 (Dual problem) とは、最適化問題における主問題 (Primary problem) の補問題であり、主問題と双対問題は表裏一体の関係にある

# ランダム化された最適化アルゴリズム

- SDCA; Stochastic Dual Coordinate Ascent

- 例えば、以下の問題  $P$ (主問題) を、次のように定める

$$\min_x c^T x \quad \text{s.t.} \quad Ax = b, x \geq 0 \quad (19)$$

- 一方の問題  $D$ (双対問題) は、次のようになる

$$\max_y b^T y \quad \text{s.t.} \quad A^T y \leq c \quad (20)$$

- 双対定理より、問題  $P$  に最適解  $x^*$  が存在すれば、問題  $D$  にも最適解  $y^*$  が存在して、 $c^T x^* = b^T y^*$  が成立

# ランダム化された最適化アルゴリズム

- SDCA; Stochastic Dual Coordinate Ascent

- 今回の場合、問題  $P$  (主問題) は、(正則化項付きの) 損失関数の最小化である
- 問題  $P$  の代わりに双対問題  $D$  を解くことができる
- このとき、問題  $P$  の解が得られるので、最適なパラメータが求まる
- SDCA は、このような考え方に基づくアルゴリズムである (詳細は省略)
- 主問題と双対問題を、交互に解くアルゴリズムも考えられる
- 以降のスライドでは、勾配  $\nabla f(w)$  の推定値に含まれるノイズを軽減するアルゴリズムをみていく

# 分散を抑えた確率的勾配降下法

- SAG; Stochastic Average Gradient

- 確率的勾配の平均を取るアルゴリズム
- SAG は、GD(勾配降下法) と SGD(確率的勾配降下法) の中間に位置
- SAG の各イテレーションでは、 $i_t \in \{1, \dots, N\}$  を選択したうえで、次の処理が実行される

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \frac{\alpha_t}{N} \sum_{i=1}^N \mathbf{y}_i^t \quad (21)$$

$$= \mathbf{w}^t - \frac{\alpha_t}{N} \left[ \sum_{i=1}^N \mathbf{y}_i^{t-1} + (\mathbf{y}_{i_t}^t - \mathbf{y}_{i_t}^{t-1}) \right] \quad (22)$$

$$= \mathbf{w}^t - \frac{\alpha_t}{N} \left[ \sum_{i=1}^N \mathbf{y}_i^{t-1} + (\nabla f_{i_t}(\mathbf{w}^t) - \mathbf{y}_{i_t}^{t-1}) \right] \quad (23)$$

# 分散を抑えた確率的勾配降下法

- 但し、 $\mathbf{y}_i^t$  は次のように定義される

$$\mathbf{y}_i^t = \begin{cases} \nabla f_i(\mathbf{w}^t) & (i = i_t) \\ \mathbf{y}_i^{t-1} & (\text{それ以外のとき}) \end{cases} \quad (24)$$

- 勾配  $\nabla f_i(\mathbf{w}^t)$  は次のように定義される

$$\nabla f_i(\mathbf{w}^t) \equiv \left. \frac{\partial}{\partial \mathbf{w}} f_i(\mathbf{w}) \right|_{\mathbf{w}=\mathbf{w}^t} \quad (25)$$

# 分散を抑えた確率的勾配降下法

- SAG; Stochastic Average Gradient

- GD では、現在のパラメータ  $w^t$  を基に、**全てのデータについて** 勾配  $\nabla f_i(w^t)$  を計算
  - ⇒ その平均  $\sum_{i=1}^N \nabla f_i(w^t)$  (**Full gradient**) を使って、パラメータ  $w$  を更新
  - ⇒ 1 回のパラメータの更新には Full gradient が必要であり、全てのデータを使って計算するため、非常に処理が重い
- SGD では、データ点  $i \in \{1, \dots, N\}$  についてのみ (**1 つのデータについてのみ**)、勾配  $\nabla f_i(w^t)$  を計算
  - ⇒ その勾配  $\nabla f_i(w^t)$  を使って、パラメータ  $w$  を更新
  - ⇒ 1 点における勾配  $\nabla f_i(w^t)$  を用いて、 $\sum_{i=1}^N \nabla f_i(w^t)$  を近似
  - ⇒ 1 回のパラメータの更新に必要な、計算量を少なく抑えられる

# 分散を抑えた確率的勾配降下法

- SAG; Stochastic Average Gradient

- SAG では、全データに対する勾配 (Full gradient) を、**少しずつ更新していく**

⇒ ある 1 つのデータ点  $i_t \in \{1, \dots, N\}$  について、現在のパラメータ  $\mathbf{w}^t$  の下で勾配  $\nabla f_{i_t}(\mathbf{w}^t)$  を計算

⇒ 新しく求めた勾配  $\nabla f_{i_t}(\mathbf{w}^t)$  を使って、以下の式で Full gradient を更新

$$\sum_{i=1}^N \mathbf{y}_i^t = \sum_{j \neq i_t} \mathbf{y}_j^{t-1} + \nabla f_{i_t}(\mathbf{w}^t) \quad (26)$$

⇒ 更新された Full gradient を使って、パラメータ  $\mathbf{w}$  を更新

- SAG では、計算される勾配にはバイアスが含まれる  
⇒ 後述の **SAGA** で計算される勾配の期待値は、勾配  $f(\mathbf{w})$  の不偏推定量に一致



# 分散を抑えた確率的勾配降下法

- SAG; Stochastic Average Gradient

- SAG では、各データにおける勾配の履歴 (現在の Full gradient)、従って  $\{\mathbf{y}_i^t\}_{i=1}^N$  を記憶しなければならないことが分かる
- 比較的小規模のニューラルネットの学習であっても、勾配を記憶するためのメモリ使用量が大きいため、アルゴリズムは使い物にならなくなる
- GD の速い収束と、SGD の速い計算時間という、双方の利点を受け継いだアルゴリズム
- SAG の改良版として、SAGA アルゴリズムが存在 (詳細は省略)
- 因みに、SAGA が何の略称なのかは不明

# 分散を抑えた確率的勾配降下法

- SAGA におけるパラメータの更新式は次のようになる [2]

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \alpha_t \left[ \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^{t-1} + (\mathbf{y}_{i_t}^t - \mathbf{y}_{i_t}^{t-1}) \right] \quad (27)$$

$$= \mathbf{w}^t - \alpha_t \left[ \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^{t-1} + (\nabla f_{i_t}(\mathbf{w}^t) - \mathbf{y}_{i_t}^{t-1}) \right] \quad (28)$$

# 分散を抑えた確率的勾配降下法

- SVRG; Stochastic Variance Reduced Gradient

- 二重ループの最適化アルゴリズムである
- 外側のループでは、Full gradient(全データについての勾配の平均)  $\nabla f(\mathbf{w}^t)$  を計算

$$\nabla f(\mathbf{w}^t) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{w}^t) = \frac{1}{N} \sum_{i=1}^N \left. \frac{\partial}{\partial \mathbf{w}} f_i(\mathbf{w}) \right|_{\mathbf{w}=\mathbf{w}^t} \quad (29)$$

- 内側のループでは、インデックス  $i \in \{1, \dots, N\}$  を選択し、次の式を用いてパラメータを更新していく

$$\mathbf{w} = \mathbf{w} - h [\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)] \quad (30)$$

- 確率的勾配  $\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t)$  は、 $\mathbf{w}$  と  $\mathbf{w}^t$  における勾配の変化  $\nabla f(\mathbf{w}) - \nabla f(\mathbf{w}^t)$  を推定するための項

# 分散を抑えた確率的勾配降下法

- SVRG が完全な勾配  $\nabla f(\mathbf{w}^t)$  を計算している間、パラメータは一度も更新されないが、同じ時間で、SGD ではパラメータが  $N$  回更新される  
⇒ 最初は、SGD の方が学習が進むことが予測される

# 分散を抑えた確率的勾配降下法

---

## Algorithm 1: SVRG; Stochastic Variance Reduced Gradient [2]

---

**parameters:**  $m$  = number of stochastic steps per epoch,  $h$  = stepsize  
(learning rate)

- 1: **for**  $s = 0, 1, \dots$  **do**
  - 2:   Compute and store full gradient  $\nabla f(\mathbf{w}^t) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{w}^t)$   
    // Full pass through data
  - 3:   Set  $\mathbf{w} = \mathbf{w}^t$
  - 4:   **for**  $t = 1$  **to**  $m$  **do**
  - 5:     Pick  $i \in \{1, \dots, N\}$ , uniformly at random
  - 6:     Update using  $\mathbf{w} = \mathbf{w} - h [\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)]$   
      // Stochastic update
  - 7:   **end for**
  - 8:    $\mathbf{w}^{t+1} = \mathbf{w}$
  - 9: **end for**
-

# 分散を抑えた確率的勾配降下法

- アルゴリズムについての補足

- SGD と比較すると SVRG の性能は良く、各イテレーションでの分散が小さい
- SVRG と、ランダム化された座標降下法 (RCD; Randomized Coordinate Descent) とを結びつけるアルゴリズムが存在
- SAGA が提唱された論文では、SAGA は SAG と SVRG の中間に位置付けられている
- SVRG、SAGA、SAG、GD を一般化したアルゴリズムが登場している

## 4 Federated Learning のための最適化アルゴリズム

# Federated Learning のための最適化アルゴリズム

- Federated Learning のための最適化アルゴリズム
  - **SVRG**(Stochastic Variance Reduced Gradient) アルゴリズムと、**DANE**(Distributed Approximate Newton) アルゴリズムを調べる
  - SVRG と DANE は一見無関係に見えるが、実は深く関連し合う
  - SVRG を改良した (Naive な)**Federated SVRG** について説明する
  - (Naive な)Federated SVRG を更に改良したアルゴリズムを、新たに提案する
    - ⇐ 各デバイスに保存された訓練データの個数、訓練データのスパース性、各デバイス上の訓練データのパターンの相違について考慮
  - Federated Learning では、訓練データが **Massively Distributed**、**Non-IID**、**Unbalanced** であるという仮定を置く
  - これに加え、**スパース性**、**Privacy Sensitive** などの仮定を設けた



# Federated Learning のための最適化アルゴリズム

- Federated Learning のためのアルゴリズムに必要な特徴

- 1 アルゴリズムの開始時に、パラメータが既に最適値であったなら、そのアルゴリズムを何度実行しても、パラメータの値が変化しない
- 2 訓練データが単一のノードにしかないとき、パラメータが収束するまでに必要な、中央のノードとの通信回数は  $\mathcal{O}(1)$  に抑えられること
- 3 データの各特徴が、単一のノードにしか現れないとき (解こうとしている問題が完全に分離され、各デバイスがパラメータの一部を学習しているとき)、 $\mathcal{O}(1)$  回の通信回数の後に、パラメータが収束すること  
  
⇐ データの各次元が、ある 1 つのノードでは 1 になるが、他の全てのノードでは 0 になるような場合
- 4 各ノードが完全に同一なデータセットを有するとき、 $\mathcal{O}(1)$  回の通信回数の後に、パラメータが収束すること

# Federated Learning のための最適化アルゴリズム

- Federated Learning のためのアルゴリズムに必要な特徴
  - 「収束する」とは、「十分に精度のある適当な解が得られる」ことを意味している
    - ⇒  $O(1)$  回とは、各デバイスと中央のノード間で、**たった 1 度だけ**やり取りすることに相当
  - (1) は、全ての最適化問題において、考慮する価値がある設定
  - (2) と (3) は、Federated Learning における極端なケース
  - (4) は、従来の最適化問題における設定 (中央の少数のノードが多量のデータを保持している状況)
    - ⇒ (4) は、Federated Learning においては最も重要でない

---

**Algorithm 2:** SVRG; Stochastic Variance Reduced Gradient (Recall) [2]

---

**parameters:**  $m$  = number of stochastic steps per epoch,  $h$  = stepsize  
(learning rate)

- 1: **for**  $s = 0, 1, \dots$  **do**
  - 2:   Compute and store full gradient  $\nabla f(\mathbf{w}^t) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{w}^t)$   
    // Full pass through data
  - 3:   Set  $\mathbf{w} = \mathbf{w}^t$
  - 4:   **for**  $t = 1$  **to**  $m$  **do**
  - 5:     Pick  $i \in \{1, \dots, N\}$ , uniformly at random
  - 6:     Update using  $\mathbf{w} = \mathbf{w} - h [\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)]$   
      // Stochastic update
  - 7:   **end for**
  - 8:    $\mathbf{w}^{t+1} = \mathbf{w}$
  - 9: **end for**
-

- SVRG; Stochastic Variance Reduced Gradient

- 外側のループでは、Full gradient(全データについての勾配の平均)  $\nabla f(\mathbf{w}^t)$  を計算
- 内側のループでは、確率的勾配によるパラメータの更新を  $m$  回実行
  - ⇒  $m$  は、データ数  $N$  の 1 倍から 5 倍程度の値に設定
  - ⇒  $m$  は、実用上は  $N$  にすることが多い
- 内側のループでは、データ点  $i$  における勾配  $\nabla f_i(\mathbf{w}), \nabla f_i(\mathbf{w}^t)$  を計算
  - ⇒ これらの勾配の差  $\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t)$  を求めて、 $\mathbf{w}$  と  $\mathbf{w}^t$  における Full gradient の差分  $\nabla f(\mathbf{w}) - \nabla f(\mathbf{w}^t)$  を推定するための項とする
  - ⇒  $\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)$  は、 $\nabla f(\mathbf{w})$  の不偏推定量を導く

- SVRG; Stochastic Variance Reduced Gradient

- 更新中の  $\mathbf{w}$  と、固定された  $\mathbf{w}^t$  との差が小さければ、  
 $\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t)$  も小さくなるので、 $\nabla f(\mathbf{w})$ (の予測値) に加わるノイズも小さくなっていると予想できる
- 内側のループを実行する度に、 $\mathbf{w}$  が  $\mathbf{w}^t$  から離れていくので、  
 $\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t)$  が増大し、従って Full gradient  $\nabla f(\mathbf{w}^t)$  に加わるノイズが増大  
⇒ 外側のループが実行され、新たな Full gradient  $\nabla f(\mathbf{w}^{t+1})$  が計算されると、Full gradient に加わるノイズは再び小さくなる
- 関数  $f = \frac{1}{N} \sum_i f_i$  が  $\lambda$ -strongly convex function で、各  $f_i$  が  $L$ -smooth function ならば、収束度合いは次のように表される (詳細は論文を参照)

$$\mathbb{E} [f(\mathbf{w}^t) - f(\mathbf{w}^*)] \leq c^t [f(\mathbf{w}^0) - f(\mathbf{w}^*)] \quad (31)$$

$\mathbf{w}^*$  は  $f(\mathbf{w})$  を最小化する最適解、 $c = \Theta\left(\frac{1}{mh}\right) + \Theta(h)$ (詳細は略)

# Federated Learning のための問題設定

- Federated Learning のための問題設定

- 解くべき問題は、経験損失  $f(\mathbf{w})$  の最小化であった

$$\min_{\mathbf{w} \in \mathbb{R}^D} f(\mathbf{w}), \quad f(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{w}) \quad (32)$$

- 各  $f_i$  は凸関数で、訓練データ  $\{\mathbf{x}_i, y_i\}_{i=1}^N$  が与えられている (多数のノードに不均等に分散)

- 分散学習のために、以下の記法を導入する

- $K$  を **ノード数** とする
- $\mathcal{P}_k$  を、ノード  $k \in \{1, \dots, K\}$  が持つ訓練データの **インデックス集合** とする
- ノード  $k$  が持つ訓練データの個数を、 $N_k = |\mathcal{P}_k|$  と表す  
 $\Rightarrow k \neq l$  のとき  $\mathcal{P}_k \cap \mathcal{P}_l = \emptyset$  (空集合)、そして  $\sum_{k=1}^K N_k = N$

- Federated Learning のための問題設定
  - 経験損失  $f(\mathbf{w})$  を、次の手順で書き直す

$$F_k(\mathbf{w}) \equiv \frac{1}{N_k} \sum_{i \in \mathcal{P}_k} f_i(\mathbf{w}) \quad (33)$$

$$\begin{aligned} f(\mathbf{w}) &= \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{w}) = \frac{1}{N} \sum_{k=1}^K \sum_{i \in \mathcal{P}_k} f_i(\mathbf{w}) \\ &= \frac{1}{N} \sum_{k=1}^K N_k \cdot \frac{1}{N_k} \sum_{i \in \mathcal{P}_k} f_i(\mathbf{w}) = \frac{1}{N} \sum_{k=1}^K N_k F_k(\mathbf{w}) \end{aligned} \quad (34)$$

- $F_k(\mathbf{w})$  は、各ノード  $k$  が最小化すべき目的関数である (凸関数)
- これより、解くべき問題は次のように書き直される

$$\min_{\mathbf{w} \in \mathbb{R}^D} f(\mathbf{w}), \quad f(\mathbf{w}) = \sum_{k=1}^K \frac{N_k}{N} F_k(\mathbf{w}) \quad (35)$$

- Federated Learning のための問題設定

- この問題を解くための最も簡単な手法は、次の通りである

$$\mathbf{w}_k^{t+1} = \arg \min_{\mathbf{w} \in \mathbb{R}^D} F_k(\mathbf{w}), \quad \mathbf{w}^{t+1} = \sum_{k=1}^K \frac{N_k}{N} \mathbf{w}_k^{t+1} \quad (36)$$

- 各ノード上で目的関数  $F_k$  を最小化し、得られた解  $\mathbf{w}_k^{t+1}$  の  $N_k$  による重み付け和を  $\mathbf{w}$  とする
- この場合、上の問題を一度だけ解けば、解  $\mathbf{w}$  が得られる ( $\mathbf{w}_k^{t+1}$  の右辺は  $t$  には依存しない) ので、各デバイスと中央のノードとの一度だけのやりとりで済む



- Federated Learning のための問題設定

- 上記のアルゴリズムは動作しない

⇐ 全体の解  $w$  が、個々の解  $w_k$  の重み付け和になっているとは考えにくい

⇐ 関数  $F_k$  の形が、全ての  $k$  について等しいならば、重み付け和にはなっている

⇐ 関数の形が全て等しいならば、単一のノード上で  $\min_{w \in \mathbb{R}^D} F_1(w)$  を解けばよいので、分散アルゴリズムを考える必要はない

- 分散アルゴリズムを導出したいが、上記のアルゴリズムでは無意味
- 但し、各ノード  $k$  が、目的関数  $F_k$  に含まれる曲がり具合の情報 (Curvature information) を最大限活用できるようにしたい

- Federated Learning のための問題設定

- 分散アルゴリズムを導出するために、各  $F_k$  に二次の項  $-(\mathbf{a}_k^t)^T \mathbf{w} + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^t\|^2$  を摂動として加算する
- そして、各ノードが次の問題を解くようにする

$$\mathbf{w}_k^{t+1} = \arg \max_{\mathbf{w} \in \mathbb{R}^D} \left( F_k(\mathbf{w}) - (\mathbf{a}_k^t)^T \mathbf{w} + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^t\|^2 \right) \quad (37)$$

$$\mathbf{w}^{t+1} = \frac{1}{K} \sum_{k=1}^K \mathbf{w}_k^{t+1} \quad (38)$$

- 各ノード  $k$  が、目的関数  $F_k$  に含まれる曲がり具合の情報 (Curvature information) を最大限活用できるようにしたい  
 $\Rightarrow$  各ノードが最適化する関数の、ヘッセ行列は  $\nabla^2 F_k + \mu \mathbf{I}$  となるので、関数  $F_k$  に含まれる勾配の情報は、殆どそのまま保存される

- Federated Learning のための問題設定

- 以下の式を解きたいが、ベクトル  $\mathbf{a}_k^t$  の決め方が分からない

$$\mathbf{w}_k^{t+1} = \arg \max_{\mathbf{w} \in \mathbb{R}^D} \left( F_k(\mathbf{w}) - (\mathbf{a}_k^t)^T \mathbf{w} + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$$

- $t \rightarrow \infty$  の極限では、 $\mathbf{w}$  が最適 ( $\mathbf{w} = \mathbf{w}^*$ ) であるとき、上式の勾配が 0 となって欲しい

$$\begin{aligned} & \nabla \left( F_k(\mathbf{w}) - (\mathbf{a}_k^t)^T \mathbf{w} + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^t\|^2 \right) \\ &= \nabla F_k(\mathbf{w}) - \mathbf{a}_k^t + \mu (\mathbf{w} - \mathbf{w}^t) = 0 \end{aligned}$$

- 即ち、 $t \rightarrow \infty$  の極限では、 $\mathbf{a}_k^t$  は次のようになって欲しい

$$\mathbf{a}_k^t = \nabla F_k(\mathbf{w}) + \mu (\mathbf{w} - \mathbf{w}^t) \simeq \nabla F_k(\mathbf{w}^*) \quad (\because \mathbf{w}^* \simeq \mathbf{w}^t)$$

- 但し、 $\mathbf{w}^*$  を知らないなので、 $\mathbf{a}_k^t = \nabla F_k(\mathbf{w}^*)$  とはできない  
 $\Rightarrow t \rightarrow \infty$  で、 $\mathbf{a}_k^t \rightarrow \nabla F_k(\mathbf{w}^*)$  となるような更新式を編み出す

# Federated Learning のための最適化アルゴリズム

- DANE; Distributed Approximate Newton
  - 先程の最適化問題は、双対問題と深く関連している
    - ⇒ 但し、上記のような問題がノード数分だけ存在するので、複雑である
  - DANE アルゴリズムでは、個々のノード上で解くための部分問題を構成することに主眼を置く
    - ⇐ 部分問題は、各ノード上のデータと、完全な勾配  $\nabla f(\mathbf{w}^t)$  にのみ依存
    - ⇐ Full gradient  $\nabla f(\mathbf{w}^t)$  は、各デバイスと中央のノード間で、1 度だけやり取りすれば計算可能
  - DANE アルゴリズムを次に示す

---

**Algorithm 3:** DANE; Distributed Approximate Newton [1]

---

**input** : regularizer  $\mu \geq 0$ , parameter  $\eta$  (default:  $\mu = 0, \eta = 1$ )

1: Initialize  $\mathbf{w}^0$

2: **for**  $t = 0, 1, \dots$  **do**

3:   Compute  $\nabla f(\mathbf{w}^t) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{w}^t)$

4:   Distribute  $\nabla f(\mathbf{w}^t)$  to all machines

5:   For each node  $k \in \{1, \dots, K\}$ , solve

$$\begin{aligned} \mathbf{w}_k^{t+1} = \arg \min_{\mathbf{w} \in \mathbb{R}^D} & \left( F_k(\mathbf{w}) - (\nabla F_k(\mathbf{w}^t) - \eta \nabla f(\mathbf{w}^t))^T \mathbf{w} \right. \\ & \left. + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^t\|^2 \right) \end{aligned} \quad (39)$$

6:   Compute  $\mathbf{w}^{t+1} = \frac{1}{K} \sum_{k=1}^K \mathbf{w}_k^{t+1}$

7: **end for**

---

- DANE; Distributed Approximate Newton

- 5 行目では、各ノードが次の部分問題を解いている

$$\mathbf{w}_k^{t+1} = \arg \min_{\mathbf{w} \in \mathbb{R}^D} \left( F_k(\mathbf{w}) - (\nabla F_k(\mathbf{w}^t) - \eta \nabla f(\mathbf{w}^t))^T \mathbf{w} + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$$

- この部分問題の解を得るためのアルゴリズムは、特に指定されていない (何でもよい)
- 問題を解くうえで、他のノードと通信する必要がない (通信コストを十分に小さくできる)
- 各ノードが、摂動を加えた最適化問題を解くアルゴリズムの 1 つ  
     $\Leftarrow \mathbf{a}_k^t$  を、 $\mathbf{a}_k^t = \nabla F_k(\mathbf{w}^t) - \eta \nabla f(\mathbf{w}^t)$  と定義している  
     $\Leftarrow \mathbf{w}^t \rightarrow \mathbf{w}^*$  ならば  $\nabla f(\mathbf{w}^t) \rightarrow \nabla f(\mathbf{w}^*) = 0$  であるので、  
     $\mathbf{a}_k^t \rightarrow \nabla F_k(\mathbf{w}^*)$  が成立

- DANE; Distributed Approximate Newton

- このアルゴリズムは、Federated Learning のために必要な条件 (2) と (3) を満たさない
- $\mu = 0, \eta = 1$  ならば、条件 (4) を満たす
- 任意の  $\mu, \eta$  について、条件 (1) を満たす
- このアルゴリズムでは、関数が 2 回微分可能であること、各ノードが独立同分布な標本を得られることを仮定
- 正則化パラメータ  $\mu$  の決め方については、改善の余地がある
  - ⇐  $\mu = 0$  であれば、ノード数  $K$  が小さいときは速やかに収束するが、 $K$  が増えるにつれて、急速に発散しやすくなる
  - ⇐  $\mu$  を大きくすれば、アルゴリズムは安定する (発散しづらくなる) が、その分パラメータの収束は遅くなる

- DANE と SVRG を融合したアルゴリズム
  - DANE アルゴリズムは、Federated Learning に適用できない (必要な条件を満たさない)
  - 部分問題の最適解を得る必要がある
    - ⇐ 簡単な問題なら可能だが、複雑な問題であれば計算コストが掛かり過ぎる
    - ⇒ 完全な最適解を得るのではなく、近似解を得るようなアルゴリズムに置き換える
    - ⇒ 部分問題を解くアルゴリズムとして、先程の **SVRG** を使用
  - SVRG アルゴリズムでは、外側のループの最初で、完全な勾配  $\nabla f(\mathbf{w}^t)$  を計算する必要があった (2 行目)
  - 完全な勾配  $\nabla f(\mathbf{w}^t)$  は、各ノードが部分問題を解く前の段階で、既に求まっている (3 行目)
    - ⇒ 各ノードでは、SVRG アルゴリズムの内側のループのみが実行され (後述)、完全な勾配を求める必要はない (既知であるとして扱う)



---

**Algorithm 4:** SVRG; Stochastic Variance Reduced Gradient (Recall) [2]

---

**parameters:**  $m$  = number of stochastic steps per epoch,  $h$  = stepsize  
(learning rate)

- 1: **for**  $s = 0, 1, \dots$  **do**
  - 2:   Compute and store full gradient  $\nabla f(\mathbf{w}^t) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{w}^t)$   
    // Full pass through data
  - 3:   Set  $\mathbf{w} = \mathbf{w}^t$
  - 4:   **for**  $t = 1$  **to**  $m$  **do**
  - 5:     Pick  $i \in \{1, \dots, N\}$ , uniformly at random
  - 6:     Update using  $\mathbf{w} = \mathbf{w} - h [\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)]$   
      // Stochastic update
  - 7:   **end for**
  - 8:    $\mathbf{w}^{t+1} = \mathbf{w}$
  - 9: **end for**
-

---

**Algorithm 5:** DANE; Distributed Approximate Newton (Recall) [1]

---

**input** : regularizer  $\mu \geq 0$ , parameter  $\eta$  (default:  $\mu = 0, \eta = 1$ )

- 1: Initialize  $\mathbf{w}^0$
- 2: **for**  $t = 0, 1, \dots$  **do**
- 3:   Compute  $\nabla f(\mathbf{w}^t) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{w}^t)$
- 4:   Distribute  $\nabla f(\mathbf{w}^t)$  to all machines
- 5:   For each node  $k \in \{1, \dots, K\}$ , solve

$$\begin{aligned} \mathbf{w}_k^{t+1} = \arg \min_{\mathbf{w} \in \mathbb{R}^D} & \left( F_k(\mathbf{w}) - (\nabla F_k(\mathbf{w}^t) - \eta \nabla f(\mathbf{w}^t))^T \mathbf{w} \right. \\ & \left. + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^t\|^2 \right) \end{aligned} \quad (40)$$

- 6:   Compute  $\mathbf{w}^{t+1} = \frac{1}{K} \sum_{k=1}^K \mathbf{w}_k^{t+1}$
  - 7: **end for**
-

- DANE と SVRG との関係性

- DANE アルゴリズムは、ある特定の条件下では、分散化されたバージョンの SVRG アルゴリズムと等価

- 以下の 2 つのアルゴリズムは等価である
- 下記の 2 つは、同一のパラメータ列  $\{w^t\}$  を出力

**1**  $\mu = 0, \eta = 1$  の下で DANE を実行し、その部分問題は SVRG アルゴリズムで解く

**2** 分散化されたバージョンの SVRG アルゴリズムを解く (後述)

- 分散化された SVRG アルゴリズム (Naive Federated SVRG) を次に示す

---

**Algorithm 6:** Naive Federated SVRG (FSVRG) [1]

---

**parameters:**  $m$  = number of stochastic steps per epoch,  $h$  = stepsize,  
data partition  $\{\mathcal{P}_k\}_{k=1}^K$

- 1: Initialize  $\mathbf{w}^0$
- 2: **for**  $t = 0, 1, \dots$ , **do**
- 3:   Compute  $\nabla f(\mathbf{w}^t) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{w}^t)$ , distribute to all machines
- 4:   **for**  $k = 1$  to  $K$  **do in parallel** over nodes  $k$  **do**
- 5:     Initialize  $\mathbf{w}_k = \mathbf{w}^t$
- 6:     **for**  $s = 1$  to  $m$  **do**
- 7:       Sample  $i \in \mathcal{P}_k$  uniformly at random
- 8:       Update using  $\mathbf{w}_k = \mathbf{w}_k - h [\nabla f_i(\mathbf{w}_k) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)]$
- 9:     **end for**
- 10:   **end for**
- 11:   Update using  $\mathbf{w}^{t+1} = \mathbf{w}^t + \frac{1}{K} \sum_{k=1}^K (\mathbf{w}_k - \mathbf{w}^t)$
- 12: **end for**

---

# Federated Learning のための最適化アルゴリズム

- 2つのアルゴリズムが等価であることの簡潔な証明
  - SVRG アルゴリズムでは、 $\nabla f(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{w})$  の不偏推定量を得るために、 $\nabla f(\mathbf{w}^t) + \nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t)$  という項を用いた
  - $\nabla f(\mathbf{w}^t)$  は、外側のループの最初で計算される (適当な時間間隔で計算され、暫くの間固定される)
  - 内側のループを実行する度に、データ  $i \in \{1, \dots, N\}$  について  $\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t)$  を計算し、 $\nabla f(\mathbf{w}^t)$  を補正する
  - DANE における部分問題を、各ノードが SVRG で解くことを考える (SVRG の **内側のループのみ** が実行される)
  - SVRG で解くのは次の問題である ( $\mu = 0, \eta = 1$ )

$$\begin{aligned} & \min_{\mathbf{w} \in \mathbb{R}^D} \left( F_k(\mathbf{w}) - (\nabla F_k(\mathbf{w}^t) - \eta \nabla f(\mathbf{w}^t))^T \mathbf{w} + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^t\|^2 \right) \\ &= \min_{\mathbf{w} \in \mathbb{R}^D} \left( F_k(\mathbf{w}) - (\nabla F_k(\mathbf{w}^t) - \nabla f(\mathbf{w}^t))^T \mathbf{w} \right) \end{aligned} \quad (41)$$

- 2つのアルゴリズムが等価であることの簡潔な証明
  - SVRG で最小化しようとしている関数の勾配を求める

$$\begin{aligned} & \frac{\partial}{\partial \mathbf{w}} \left( F_k(\mathbf{w}) - (\nabla F_k(\mathbf{w}^t) - \nabla f(\mathbf{w}^t))^T \mathbf{w} \right) \\ &= \nabla F_k(\mathbf{w}) - \nabla F_k(\mathbf{w}^t) - \nabla f(\mathbf{w}^t) \end{aligned} \quad (42)$$

- 上式について、 $\nabla f(\mathbf{w}^t)$  は、DANE アルゴリズムの 3 行目で既に求まっている
- 残りの  $\nabla F_k(\mathbf{w})$  と  $\nabla F_k(\mathbf{w}^t)$  の計算について考える
- $\nabla F_k(\mathbf{w}) = \sum_{i \in \mathcal{P}_k} \nabla f_i(\mathbf{w})$  であるから、 $\nabla F_k(\mathbf{w})$  は、 $i \in \mathcal{P}_k$  である全ての  $i$  についての、勾配  $\nabla f_i(\mathbf{w})$  の足し合わせである ( $\nabla F_k(\mathbf{w}^t)$  についても同様)
- SVRG の内側のループでは、ある 1 つのインデックス  $i \in \mathcal{P}_k$  をランダムに選び出し、パラメータ  $\mathbf{w}$  を更新する

# Federated Learning のための最適化アルゴリズム

- 2つのアルゴリズムが等価であることの簡潔な証明
  - SVRG の内側のループでは、ある 1つのインデックス  $i \in \mathcal{P}_k$  をランダムに選び出し、パラメータ  $\mathbf{w}$  を更新する

⇒ 内側のループでは、1つのデータ点  $i$  についての確率的勾配だけを使って、 $\nabla F_k(\mathbf{w}) - \nabla F_k(\mathbf{w}^t) - \nabla f(\mathbf{w}^t)$  を推定したい

⇒  $\nabla f(\mathbf{w}^t)$  は既知であるから、残りの項  $\nabla F_k(\mathbf{w}) - \nabla F_k(\mathbf{w}^t)$  を、確率的勾配  $\nabla f_i(\mathbf{w}), \nabla f_i(\mathbf{w}^t)$  を使って近似する

- SVRG のときと同じ方法で、次のように近似できる

$$\nabla F_k(\mathbf{w}) = \nabla F_k(\mathbf{w}^t) + (\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t)) \quad (43)$$

$$\nabla F_k(\mathbf{w}^t) = \nabla F_k(\mathbf{w}^t) + (\nabla f_i(\mathbf{w}^t) - \nabla f_i(\mathbf{w}^t)) \quad (44)$$

# Federated Learning のための最適化アルゴリズム

- 2つのアルゴリズムが等価であることの簡潔な証明
  - 以下は単に、 $f(\mathbf{w})$  を  $F_k(\mathbf{w})$  に置き換えているだけ

$$\nabla F_k(\mathbf{w}) = \nabla F_k(\mathbf{w}^t) + (\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t)) \quad (45)$$

- $\nabla F_k(\mathbf{w}^t)$  を、内側のループが実行される度に計算するのは、計算コストの観点から避けたいので、適当な時間間隔で求めて、暫くの間は使い回すことにする

⇒ 現在のパラメータ  $\mathbf{w}$  についての勾配  $\nabla F_k(\mathbf{w})$  の不偏推定量を得るために、古いパラメータ  $\mathbf{w}^t$  についての勾配  $\nabla F_k(\mathbf{w}^t)$  を、データ  $i$  についての項  $\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t)$  で補正する



- 2つのアルゴリズムが等価であることの簡潔な証明

- このような近似の下で、内側のループで計算される勾配は、次のようになる

$$\nabla F_k(\mathbf{w}) - \nabla F_k(\mathbf{w}^t) - \nabla f(\mathbf{w}^t) \quad (46)$$

$$\begin{aligned} &= [\nabla F_k(\mathbf{w}^t) + \nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t)] \\ &\quad - [\nabla F_k(\mathbf{w}^t) + \nabla f_i(\mathbf{w}^t) - \nabla f_i(\mathbf{w}^t)] - \nabla f(\mathbf{w}^t) \end{aligned} \quad (47)$$

$$= \nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t) \quad (48)$$

- これより、DANE アルゴリズムの部分問題に、SVRG アルゴリズムを適用したとき、各ノードにおけるパラメータの更新式は次のようになる

$$\mathbf{w} = \mathbf{w} - h [\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)] \quad (49)$$

- 上の更新式は、FSVRG(Algorithm 6) の 8 行目と等しい  
⇒ DANE( $\mu = 0, \eta = 1$ ) の部分問題に SVRG を組み込んだアルゴリズムは、(Naive な)FSVRG と等価

- Naive Federated SVRG についての補足

- 外側のループについて、最後の更新式が次のようになっている

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \frac{1}{K} \sum_{k=1}^K (\mathbf{w}_k - \mathbf{w}^t) \quad (50)$$

- これは、DANE における次の式に等しい

$$\mathbf{w}^{t+1} = \frac{1}{K} \sum_{k=1}^K \mathbf{w}_k^{t+1} \quad (51)$$

- $(\mathbf{w}_k - \mathbf{w}^t)$  は、各ノードから中央のサーバに送られる、パラメータの差分を表している  
⇒  $\mathbf{w}^t$  が中央のサーバから送られてきたパラメータ、 $\mathbf{w}_k$  はノード  $k$  上の訓練データを使って更新されたパラメータ  
⇒  $\mathbf{w}_k$  をそのまま送っても良いが、差分を送った方がデータを圧縮できる (通信量を削減)

- SVRG アルゴリズムの導出

- (Naive な)Federated SVRG を更に改良したアルゴリズムを、ここでは提案する
  - ⇐ 各デバイスに保存された訓練データの個数、訓練データのスパース性、データが独立同分布でない (Non-IID) ことについて考慮
- 各デバイス間では、訓練データの個数が大きく異なる
- データのある特徴は、ごく少数のノードにだけ出現し、それ以外の大多数のノードには出現しないかもしれない
- 各デバイス上のデータは、データ全体の分布を表していない (データ全体は、デバイスの地理的な位置やタイムゾーンなどが影響し、特定のパターンによってクラスタ化されているかもしれない)

- SVRG アルゴリズムの導出 (記法の整理)

- $N$ : 訓練データの個数,  $K$ : ノード数,  $D$ : パラメータの次元
- $\mathcal{P}_k$ : ノード  $k$  が保持する訓練データのインデックスの集合  
 $\Rightarrow \mathcal{P}_k \subseteq \{1, \dots, N\}, \forall k \neq l \mathcal{P}_k \cap \mathcal{P}_l = \emptyset$
- $N_k = |\mathcal{P}_k|$ : ノード  $k$  が保持する訓練データの個数  
 $\Rightarrow \sum_{k=1}^K N_k = N$
- $N^j = |\{i \in \{1, \dots, N\} | \mathbf{x}_i^T \mathbf{e}_j \neq 0\}|$ :  $j$  番目の次元が 0 ではない訓練データの個数
- $N_k^j = |\{i \in \mathcal{P}_k | \mathbf{x}_i^T \mathbf{e}_j \neq 0\}|$ : ノード  $k$  が保持する、 $j$  番目の次元が 0 ではない、訓練データの個数
- $\phi^j = N^j/N$ : 全訓練データのうち、 $j$  番目の次元が 0 でないものの割合
- $\phi_k^j = N_k^j/N_k$ : ノード  $k$  が保持する全訓練データのうち、 $j$  番目の次元が 0 でないものの割合

# Federated Learning のための最適化アルゴリズム

- $s_k^j = \phi^j / \phi_k^j$ : 全訓練データと、ノード  $k$  が保持する訓練データにおける、 $j$  番目の次元が 0 でない割合の比率  
⇒  $s_k^j$  が小さいとき、ノード  $k$  は、次元  $j$  に関する訓練データを、他のノードよりも多く持っていることを意味している  
⇒  $s_k^j$  の逆数は、ノード  $k$  が持っている訓練データは、次元  $j$  に関してどの程度レアかを表すと考えられる
- $S_k = \text{diag}(s_k^j)$ :  $s_k^j$  を並べた対角行列 ( $S_k \in \mathbb{R}^{D \times D}$ )
- $\omega^j = \left| \left\{ \mathcal{P}_k | n_k^j \neq 0 \right\} \right|$ : 次元  $j$  が 0 ではない訓練データを持っている、ノードの数
- $a^j = K / \omega^j$ : 次元  $j$  が 0 でない訓練データが、ノードに出現する割合  
⇒ 次元  $j$  のレア度を表す
- $A = \text{diag}(a^j)$ :  $a^j$  を並べた対角行列 ( $A \in \mathbb{R}^{D \times D}$ )
- 提案された Federated SVRG アルゴリズムを、次に示す

---

**Algorithm 7: Federated SVRG (FSVRG) [1]**

---

**parameters:**  $h$  = stepsize, data partition  $\{\mathcal{P}_k\}_{k=1}^K$ , diagonal matrix  $\mathbf{A}, \mathbf{S}_k \in \mathbb{R}^{D \times D}$  for  $k \in \{1, \dots, K\}$

```
1: for  $t = 0, 1, \dots$  do
2:   Compute  $\nabla f(\mathbf{w}^t) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{w}^t)$ 
3:   for  $k = 1$  to  $K$  do in parallel over nodes  $k$  do
4:     Initialize  $\mathbf{w}_k = \mathbf{w}^t$ ,  $h_k = h/N_k$ 
5:     Let  $\{i_s\}_{s=1}^{N_k}$  be random permutation of  $\mathcal{P}_k$ 
6:     for  $s = 1, \dots, N_k$  do
7:        $\mathbf{w}_k = \mathbf{w}_k - \textcolor{red}{h}_k [\textcolor{red}{S}_k [\nabla f_{i_s}(\mathbf{w}_k) - \nabla f_{i_s}(\mathbf{w}^t)] + \nabla f(\mathbf{w}^t)]$ 
8:     end for
9:   end for
10:   $\mathbf{w}^t = \mathbf{w}^t + \textcolor{red}{A} \sum_{k=1}^K \frac{N_k}{N} (\mathbf{w}_k - \mathbf{w}^t)$ 
11: end for
```

---

# Federated Learning のための最適化アルゴリズム

- Federated SVRG アルゴリズムの改良点

- 1 ノードごとにステップサイズを変更:  $h_k = h/N_k$
  - 2 各ノードが保持するデータ数に比例した更新量:  $\frac{N_k}{N} (\mathbf{w}_k - \mathbf{w}^t)$
  - 3 確率的勾配の各次元に対するスケーリング:  $\mathbf{S}_k$
  - 4 パラメータの差分の各次元に対するスケーリング:  $\mathbf{A} (\mathbf{w}_k - \mathbf{w}^t)$
- 個人的に、このようなアルゴリズムの細工にはあまり魅力を感じない
  - このような細工よりも、単純にデータ数を増やせば良いんじゃないか?
  - このアルゴリズムは、他のアルゴリズムよりも圧倒的に収束が速い

# Federated Learning のための最適化アルゴリズム

## 1 ノードごとにステップサイズを変更: $h_k = h/N_k$

- Naive Federated SVRG では、各ノードでのパラメータの更新回数は、 $m$  回に統一されていた  
⇒ しかし、各ノードが保持する訓練データの数大きく異なるのだから、パラメータの更新回数を、全ノードにわたって同じにするのは良くない  
⇒ それゆえ、各ノードでは、自身が持つ全ての訓練データを使って、パラメータを更新
- ステップサイズを  $N_k$  に反比例するように設定し、各ノードのパラメータ更新量が同程度の大きさになるようにする



# Federated Learning のための最適化アルゴリズム

- 2 各ノードが保持するデータ数に比例した更新量:  $\frac{N_k}{N} (\mathbf{w}_k - \mathbf{w}^t)$
- Naive Federated SVRG では、以下のようにパラメータ  $\mathbf{w}$  が更新された

$$\begin{aligned}\mathbf{w}^{t+1} &= \mathbf{w}^t + \frac{1}{K} \sum_{k=1}^K (\mathbf{w}_k - \mathbf{w}^t) \\ &= \mathbf{w}^t + \frac{h}{K} \sum_{k=1}^K \sum_i [\nabla f_i(\mathbf{w}_k) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)] \\ &= \mathbf{w}^t + \frac{h}{K} \sum_{k=1}^K \sum_i [\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)] \\ &\simeq \mathbf{w}^t + \frac{h}{K} \sum_{k=1}^K \mathbb{E} [\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)] \\ &= \mathbf{w}^t + \frac{h}{K} \mathbb{E} \left[ \sum_{k=1}^K (\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)) \right]\end{aligned}$$

# Federated Learning のための最適化アルゴリズム

- 上式における  $i$  は、 $\mathcal{P}_k$  から重複を許して、適当に選択された  $m$  個のインデックスであるとする (但し、最後の式における  $i$  は、 $\mathcal{P}_k$  から適当に選択されたインデックスである)
- また、ある時点において、全ての  $k \in \{1, \dots, K\}$  について  $\mathbf{w}_k = \mathbf{w}$  であるとしている
- ここでは、適当な  $\alpha_k$  を導入して、 $\nabla f(\mathbf{w}^t)$  の不偏推定量が得られるようにしたい ( $i$  は  $\mathcal{P}_k$  から適当に選択されたもの)  
     $\Leftarrow$  勾配降下法などの、目的関数の勾配のみを用いたアルゴリズム (一次法; First-order method) に対して求められる性質
- 即ち、以下が成り立って欲しい ( $\mathbf{w}^t$  に足される量が、 $\nabla f(\mathbf{w}^t)$  に比例する)

$$\mathbb{E} \left[ \sum_{k=1}^K \alpha_k (\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)) \right] = \nabla f(\mathbf{w}^t) \quad (52)$$

# Federated Learning のための最適化アルゴリズム

- $\alpha_k$  は次のようにして求められる

$$\begin{aligned} & \mathbb{E} \left[ \sum_{k=1}^K \alpha_k (\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)) \right] \\ &= \sum_{k=1}^K \alpha_k \mathbb{E} [\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)] \\ &\simeq \sum_{k=1}^K \alpha_k \frac{1}{N_k} \sum_{i \in \mathcal{P}_k} [\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)] \quad (53) \end{aligned}$$

- $\alpha_k = N_k/N$  とすることで以下を得る

$$= \frac{1}{N} \sum_{k=1}^K \sum_{i \in \mathcal{P}_k} [\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)] \simeq \nabla f(\mathbf{w})$$

# Federated Learning のための最適化アルゴリズム

- これより、各ノードで計算されたパラメータの差分が、 $N_k$  に比例するように、更新量を調節することが考えられる

$$\begin{aligned}\mathbf{w}^{t+1} &= \mathbf{w}^t + \frac{h}{K} \mathbb{E} \left[ \sum_{k=1}^K \frac{N_k}{N} (\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)) \right] \\ &= \mathbf{w}^t + \frac{h}{K} \frac{N_k}{N} \mathbb{E} \left[ \sum_{k=1}^K (\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)) \right]\end{aligned}$$

## 3 確率的勾配の各次元に対するスケーリング: $S_k$

- 行列  $S_k$  は、 $s_k^j = \phi^j / \phi_k^j$  を対角成分にもつ
- データ数が  $N = 10^6$ 、ノード数が  $K = 10^3$  であるとする
- 次元  $j$  が非零になるデータは  $10^3$  個あるが、その全てが、ある 1 つのノード  $k$  に集中しているとする
- このとき  $s_k^j = \phi^j / \phi_k^j = 10^{-3}$  であるので、ノード  $k$  でパラメータ  $w$  が更新されるとき、次元  $j$  の更新量は  $10^3$  分の 1 倍される
- 行列  $S_k$  でスケーリングを行わない場合、ノード  $k$  における、次元  $j$  に対する更新量は、他のノードの  $10^3$  倍程度になる (或いは、パラメータ  $w$  の次元  $j$  は、ノード  $k$  でしか更新されない)  
⇒ パラメータ  $w$  の次元  $j$  について、急速に発散してしまう恐れがある
- パラメータ  $w$  の更新に使用する、勾配の (推定量の) 各次元が、大体同じになるように揃える役割 (個人的には、まだよく理解できていない)

# Federated Learning のための最適化アルゴリズム

## 4 パラメータの差分の各次元に対するスケーリング: $A(w_k - w^t)$

- 行列  $A$  は、 $a_j = K/\omega_j$  を対角成分にもつ
- 行列  $A$  の各成分は、 $\omega_j$  に反比例している  
⇒ 次元  $j$  が非零であるデータを持つノード数  $\omega_j$  が少ないほど、パラメータ  $w$  の次元  $j$  に対する更新量は大きくなる

⇒  $\omega_j$  が小さいとき、少ないノードから、次元  $j$  に関する情報が伝達されているので、その情報の価値は高い (パラメータ  $w$  の次元  $j$  の更新に際して、大いに役立つ情報である) と考えられる (従って更新量を大きくする)

- 4つの改善のうち、どれが最も効果的であったかは不明であるほか、これらの改善策が、相互に干渉し合うと思われる
- $h_k, \frac{N_k}{N}, S_k, A$  によって、適切なスケーリングを行うことで、各ノード上の訓練データ数のばらつき、データのスパース性などに対処できそうなのは、直感的には分かる

## 5 最適化アルゴリズムの比較

# 最適化アルゴリズムの比較

## ● 最適化アルゴリズムの比較

- Google+の投稿に、1 つ以上のコメントが付くかどうかを予測する、二値分類のタスク
- L2 正則化項を含めたロジスティック回帰
- Google+の投稿から作成したデータセットを、実験に使用
- 各ユーザの投稿が、各ユーザの使うデバイス上に保存されているという状況を想定
- $10^2$  以上のパブリックな投稿を英語で行っている、 $10^4$  のユーザをランダムに選択
- 選択されたユーザの投稿を集めて、そのうちの 75% にあたる  $N = 2,166,693$  個を訓練データとした
- 入力データとなる投稿は Bag-of-words 形式に変換された  
⇐ (全投稿データを基に得られた) 最頻出単語 20,000 語と、それ以外の不明な単語の出現回数と、バイアス項を加えたことで、パラメータ数は  $D = 20,002$  となった



# 最適化アルゴリズムの比較

## ● 最適化アルゴリズムの比較

- 各ユーザが投稿する内容は一般に異なるので、各ユーザが保持するデータの特徴も全く異なる
  - ⇒ 各ユーザの投稿内容は、全データの分布内で、クラスタを形成していると考えられる
  - ⇒ 各ユーザの保持するデータは、独立同分布標本とは**仮定できない**
  - ⇒ もしそのような仮定をすると、全てのユーザが、あらゆる内容の投稿を満遍なくしていることになる
- 入力データは Bag-of-words 形式であるから、非常にスパースである
  - ⇒ 殆どの投稿は、ごく一部の単語しか含んでいないので、入力データの中で非零になる要素数が少ない
  - ⇒ 殆どの次元は、非零になる (その次元に対応する特徴が出現する) 確率が低い
  - ⇒ 入力データの各特徴の出現確率を、次の図 1 に示す (不明な単語を表す次元の出現率が高くなっているが、これは実際とは異なる)

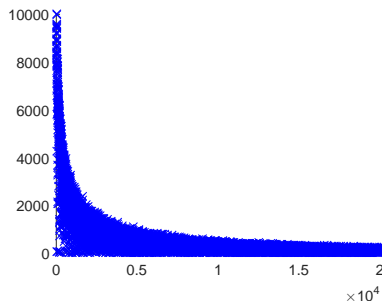


Figure 1: Features vs. appearance on nodes. The  $x$ -axis is a feature index, and the  $y$ -axis represents the number of nodes where a given feature is present.

図 1: 特徴の出現確率のヒストグラム [1]

# 最適化アルゴリズムの比較

## ● 最適化アルゴリズムの比較

- Google+の投稿に、1 つ以上のコメントが付くかどうかを予測する、二値分類のタスク
- L2 正則化項を含めたロジスティック回帰の、テストデータに対する誤差 (分類誤り) を示す
- 全ての投稿に対して、コメントが付かない ( $-1$ ) という予測をすると、誤差は 33.16%
- 入力データを全て用いて、通常のロジスティック回帰を行うと、誤差は 26.27%
- 各ユーザの全ての投稿に対して、同一の予測をすると、誤差は 17.14%  
⇒ 個々の投稿内容ではなく、誰が投稿しているのかに着目して、コメントが付くかどうかの予測を行った方が、精度が良くなることを示唆 (直感的にもそうである)

# 最適化アルゴリズムの比較

- 最適化アルゴリズムの比較

- 全ユーザで共通のモデルを基に、各ユーザごとにチューニングすることで、精度を向上させられる可能性  
⇒ もし精度が大幅に向上するのであれば、各ノードごとに、異なるデータのパターンが出現していることも確認できる
- 各最適化アルゴリズムでの性能比較を行う
- 提案手法 (Federated SVRG; Algorithm 7) のハイパーパラメータは、ステップサイズ  $h$  のみであるが、性能が最も良くなるものを選んだ

# 最適化アルゴリズムの比較

## ● 最適化アルゴリズムの比較

- オフラインアルゴリズムで得られる最大の性能 (OPT)、勾配降下法 (GD)、CoCoA+アルゴリズム (詳細は略)、Federated SVRG アルゴリズム (FSVRG)、ランダムにデータをシャッフルした上で行う Federated SVRG (FSVRGR) の性能を比較
  - GD: 個々のノードが、自身が持つ訓練データを全て使って、勾配降下法を行っている
  - CoCoA+: このスライドでは省略
  - FSVRG: この論文における提案手法
  - FSVRGR: 各ノードが保持するデータ数は変化させない (ばらつきは維持する) が、データを一旦全て集めて、ランダムにシャッフルした上で、個々のノードに戻している (各ノードが保持するデータが IID ではなくても、FSVRG アルゴリズムが動作することを示すために用意してある)
- 目的関数の推移、またはテストデータに対する分類誤差の推移を縦軸に、そして各デバイスと中央のサーバとの一連のやり取り (Round) の回数を横軸に取ったグラフを、次の図 2 に示す

# 最適化アルゴリズムの比較

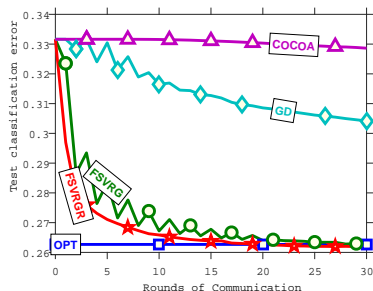
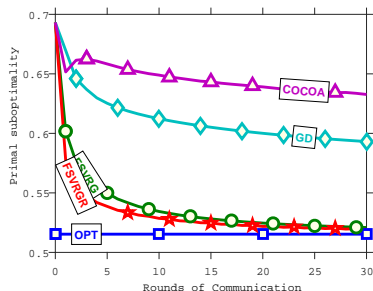


Figure 2: Rounds of communication vs. objective function (left) and test prediction error (right).

図 2: 最適化アルゴリズムの性能比較 [1]

# 最適化アルゴリズムの比較

- 最適化アルゴリズムの比較

- FSVRG は、僅か 30 回程度の Round で収束しており、Federated Optimization における課題を最初に解決したアルゴリズムといえる

⇒ CoCoA+や、それ以外の通信効率の良い (Communication-Efficient) 分散アルゴリズムでも、収束していない (学習が上手く行っていない)

⇒ 学習が上手く行かないのは、他のアルゴリズムでは、全てのデバイスが独立同分布なデータをもつと仮定しているため

- FSVRG と FSVRGR の結果の差はごく僅かであり、Naive Federated SVRG を改良することで、各デバイスが持つデータの分布に左右されない (ロバストな) アルゴリズムを得られた

## 6 結論



- Federated Learning の特徴
  - Massively Distributed、Non-IID、Unbalanced、Privacy Sensitive、Sparse の 5 つが挙げられる
  - これらの性質をもつ厳しい環境下においても、効率的に学習が進むアルゴリズムを設計可能
- 今後の研究の方向性について
  - 非同期版のアルゴリズムの開発や、アルゴリズムの理論的な解明 (特に収束性)
  - 非凸関数の最適化 (ニューラルネットが代表例) についての理論的な解明
  - 全ユーザ間で共有されるモデルと、個々のユーザに特化したモデル双方の活用

- [1] Jakub Konečný, H. Brendan McMahan, Daniel Ramage, and Peter Richtárik.  
Federated optimization: Distributed machine learning for on-device intelligence.  
[arXiv:1610.02527](https://arxiv.org/abs/1610.02527), 2016.
- [2] Pradeep Ravikumar.  
Stochastic optimization methods.  
[http://www.cs.cmu.edu/~pradeepr/convexopt/Lecture\\_Slides/stochastic\\_optimization\\_methods.pdf](http://www.cs.cmu.edu/~pradeepr/convexopt/Lecture_Slides/stochastic_optimization_methods.pdf), 2017.
- [3] Suzuki Taiji.  
機械学習におけるオンライン確率的最適化の理論.  
<https://www.slideshare.net/trinmu/stochasticoptim2013>, 2013.