

# Intelligence Artificielle

## Recherche Opérationnelle

Jérôme MALICK

CNRS, Lab. Jean Kuntzmann & MIAI (Institut IA de Grenoble)



Tutoriel du GdR RO

ROADEF – Feb. 2020 – Montpellier

Fact: AI is everywhere...

On the news

...and people's discussions



Fact: AI is everywhere...

On the news

...and people's discussions

In our universities

e.g. AI institutes



Fact: AI is everywhere...

On the news

...and people's discussions



In our universities

e.g. AI institutes



In the scientific community

e.g. Turing award 2019

to deep learning pioneers





Let's play together !

Various, fruitful interactions AI  $\Leftrightarrow$  OR

---

Let's play together !

Various, fruitful interactions AI  $\Leftrightarrow$  OR

- E.g. last year ROADEF/EURO Challenge



Let's play together !

Various, fruitful interactions AI  $\Leftrightarrow$  OR

- E.g. last year ROADEF/EURO Challenge



New perspectives in decision-making 😊

- E.g. design and control of smart infrastructure (grids, transportation,...)
- Hence the title of this “tutorial” (...actually more like a bird-eye overview)

## Vast domain at the interface AI/OR

A couple of pointers in a jungle of references



A. Parmentier's talk



B. Rottembourg

## Vast domain at the interface AI/OR

A couple of pointers in a jungle of references



A. Parmentier's talk



B. Rottembourg

A complementary viewpoint:

This tutorial : modest goals

- recall basics of machine learning
- mention (important) general ideas
- illustrate related theoretical research topics

## Recent research in the team DAO @LJK



Y.-G. Hsieh, F. Iutzeler, J. Malick, P. Mertikopoulos,  
On the Convergence of Single-Call Stochastic Extra-Gradient Methods  
NeurIPS, 2019



M. Grishchenko, F. Iutzeler, J. Malick  
Subspace Descent Methods with Identification-Adapted Sampling  
Submitted to: Maths of OR, 2019



Y. Laguel, J. Malick, Z. Harchaoui  
Superquantile Minimization: Oracles and First-order Algorithms  
Submitted to: Optimization Methods and Software, 2019

# Outline

- ➊ Back to basics: learning is optimizing
- ➋ Discussion: (some) perspectives, (deep) questions, and (personal) thoughts
- ➌ Highlight: fresh interest on min-max
- ➍ Highlight: collaborative/federative learning

# Outline

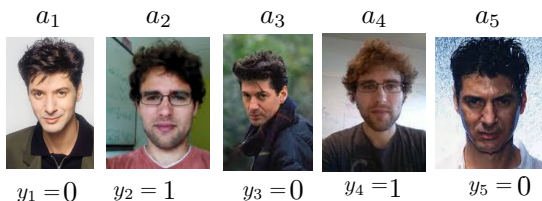
- 1 Back to basics: learning is optimizing
- 2 Discussion: (some) perspectives, (deep) questions, and (personal) thoughts
- 3 Highlight: fresh interest on min-max
- 4 Highlight: collaborative/federative learning



## Supervised learning set-up

**Data:**  $n$  observations  $(a_i, y_i) \in \mathbb{R}^m \times \mathcal{Y}$

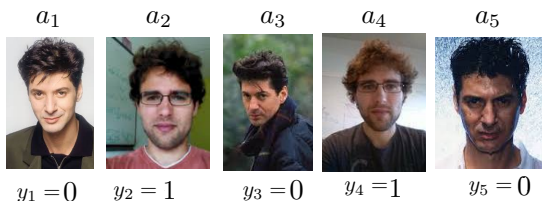
**Task:** e.g. binary classification ( $\mathcal{Y} = \{0, 1\}$ )



## Supervised learning set-up

**Data:**  $n$  observations  $(a_i, y_i) \in \mathbb{R}^m \times \mathcal{Y}$

**Task:** e.g. binary classification ( $\mathcal{Y} = \{0, 1\}$ )



**Model:** for a new  $a$ , prediction  $h(a, x) \in \mathcal{Y}$  parameterized by  $x \in \mathbb{R}^d$

usually  $x = \beta$  (in stats)  $x = \omega$  (in learning) or  $x = \theta$  (in deep learning)

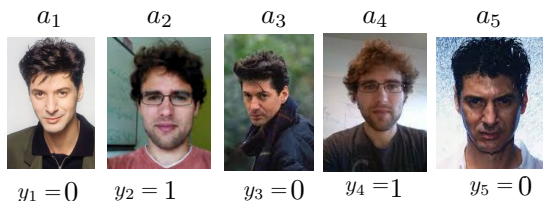
**Standard prediction functions :**

- Linear prediction:  $h(a, x) = a^\top x$

# Supervised learning set-up

**Data:**  $n$  observations  $(a_i, y_i) \in \mathbb{R}^m \times \mathcal{Y}$

**Task:** e.g. binary classification ( $\mathcal{Y} = \{0, 1\}$ )



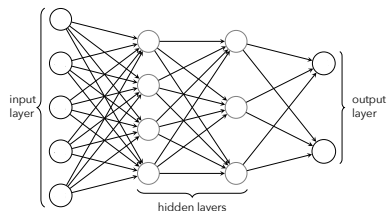
**Model:** for a new  $a$ , prediction  $h(a, x) \in \mathcal{Y}$  parameterized by  $x \in \mathbb{R}^d$

usually  $x = \beta$  (in stats)  $x = \omega$  (in learning) or  $x = \theta$  (in deep learning)

## Standard prediction functions :

- Linear prediction:  $h(a, x) = a^\top x$
- Artificial neural networks:  

$$h(a, x) = x_m^\top \sigma(x_{m-1}^\top \cdots \sigma(x_1^\top a))$$



## Optimization comes into play

- Learning = finding the “best” parameter  $\bar{x}$ 
  - = finding  $\bar{x}$  such that  $h(a_i, \bar{x}) \simeq y_i$  (and generalizes well on unseen data)
  - = solving an optimization problem
- Regularized empirical risk minimization

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(a_i, x))$$

## Optimization comes into play

- Learning = finding the “best” parameter  $\bar{x}$ 
  - = finding  $\bar{x}$  such that  $h(a_i, \bar{x}) \simeq y_i$  (and generalizes well on unseen data)
  - = solving an optimization problem
- Regularized empirical risk minimization  
(regularization avoids overfitting, helps numerically, or imposes structure to  $x$ )

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(a_i, x)) + \lambda r(x)$$

# Optimization comes into play

- Learning = finding the “best” parameter  $\bar{x}$   
 = finding  $\bar{x}$  such that  $h(a_i, \bar{x}) \simeq y_i$  (and generalizes well on unseen data)  
 = solving an optimization problem

- **Regularized empirical risk minimization**

(regularization avoids overfitting, helps numerically, or imposes structure to  $x$ )

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(a_i, x)) + \lambda r(x)$$

- Example: linear model  $h(a, x) = a^\top x$  and **least-squares loss**

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (y_i - a_i^\top x)^2 + \frac{\lambda}{2} \|x\|_2^2 = \frac{1}{2n} \|Ax - y\|_2^2 + \frac{\lambda}{2} \|x\|_2^2$$

# Optimization comes into play

- Learning = finding the “best” parameter  $\bar{x}$ 
  - = finding  $\bar{x}$  such that  $h(a_i, \bar{x}) \simeq y_i$  (and generalizes well on unseen data)
  - = solving an optimization problem
- **Regularized empirical risk minimization**

(regularization avoids overfitting, helps numerically, or imposes structure to  $x$ )

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(a_i, x)) + \lambda r(x)$$

- Example: linear model  $h(a, x) = a^\top x$  and **least-squares loss**

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (y_i - a_i^\top x)^2 + \frac{\lambda}{2} \|x\|_2^2 = \frac{1}{2n} \|Ax - y\|_2^2 + \frac{\lambda}{2} \|x\|_2^2$$

$$\min_{x \in \mathbb{R}^d} \frac{1}{2n} \|Ax - y\|_2^2 + \frac{\lambda}{2} \|x\|_1 \quad \ell_1\text{-norm promotes sparse solutions}$$

## Stochastic gradient rules

(in the simple case  $r = 0$ )

- Basic optimization algorithm : stochastic gradient descent (SGD)

Draw random  $i_k$

$$x_{k+1} = x_k - \gamma_k g_k \quad \text{with} \quad g_k = \nabla \ell(y_{i_k}, h(a_{i_k}, x_k))$$

$$\text{with} \quad \mathbb{E}[g_k] = \nabla f(x_k)$$



# Stochastic gradient rules

(in the simple case  $r = 0$ )

- Basic optimization algorithm : stochastic gradient descent (SGD)

Draw random  $i_k$

$$x_{k+1} = x_k - \gamma_k g_k \quad \text{with} \quad g_k = \nabla \ell(y_{i_k}, h(a_{i_k}, x_k))$$

$$\text{with} \quad \mathbb{E}[g_k] = \nabla f(x_k)$$

- We can often compute the gradient

E.g. back-propagation for neural networks to derivate  $x \mapsto \ell(y, h(a, x))$

- Tuning of  $\gamma_k$  is the key of efficiency
- Zoology: SAG, SDCA, Miso, SVRG, SAGA, Adam, AdaGrad, Eve,...  
+ mini-batch + prox-versions + accelerated versions + 2nd order

## Bottomline

Stochastic first-order optimization (training) methods work great!

## Stochastic gradient rules

(in the simple case  $r = 0$ )

- Basic optimization algorithm : stochastic gradient descent (SGD)

Draw random  $i_k$

$$x_{k+1} = x_k - \gamma_k g_k \quad \text{with} \quad g_k = \nabla \ell(y_{i_k}, h(a_{i_k}, x_k))$$

$$\text{with} \quad \mathbb{E}[g_k] = \nabla f(x_k)$$

- We can often compute the gradient

E.g. back-propagation for neural networks to derivate  $x \mapsto \ell(y, h(a, x))$

- Tuning of  $\gamma_k$  is the key of efficiency
- Zoology: SAG, SDCA, Miso, SVRG, SAGA, Adam, AdaGrad, Eve,...  
+ mini-batch + prox-versions + accelerated versions + 2nd order

### Bottomline

Stochastic first-order optimization (training) methods work great!

- ...except when they don't
- non-convex landscape, very few guarantees...
  - leaving aside i.i.d. train/test data, biased data,...

# More than the just workhorse...

## Optimization plays a fundamental role in learning

### E.g. Test of Time Award

NeurIPS 2018

NeurIPS 2019

ICML 2019

#### The Tradeoffs of Large Scale Learning

**Leon Bottou**  
NEC Laboratories of America  
Princeton, NJ 08540, USA  
leonbottou.org

**Olivier Bousquet**  
Google Zurich  
8002 Zurich, Switzerland  
olivier.bousquet@gmx.org

#### Abstract

This contribution develops a theoretical framework that takes into account the effect of approximate optimization on learning algorithms. The analysis shows distinct tradeoffs for the case of small-scale and large-scale learning problems. Small-scale learning problems are subject to the usual approximation-optimization tradeoff. Large-scale learning problems are subject to a qualitatively different tradeoff involving the computational complexity of the underlying optimization algorithms in non-trivial ways.

#### 1 Motivation

The computational complexity of learning algorithms has seldom been taken into account by the learning theory. Vapnik [1] states that a problem is “learnable” when there exists a probably approximately correct learning algorithm with polynomial complexity. Whereas much progress has been made on the statistical aspect (e.g., [2, 3, 4]), very little has been told about the complexity side of this proposal (e.g., [5]).

Computational complexity becomes the limiting factor when one envisions large amounts of training data. Two important examples come to mind:

#### Dual Averaging Method for Regularized Stochastic Learning and Online Optimization

**Lin Xiao**  
Microsoft Research, Redmond, WA 98052  
lin.xiao@microsoft.com

#### Abstract

We consider regularized stochastic learning and online optimization problems, where the objective function is the sum of two convex terms: one is the loss function of the learning task, and the other is a simple regularization term such as norm for promoting sparsity. We develop a new online algorithm, the regularized dual averaging (RDA) method, that can explicitly exploit the regularization structure in an online setting. In particular, at each iteration, the learning variables are adjusted by solving a simple optimization problem that involves the minimization of all past subgradients of the loss functions and the whole regularization term, not just its subgradient. Computational experiments show that the RDA method can be very effective for sparse online learning with regularization.

#### 1 Introduction

In machine learning, online algorithms operate by repeatedly drawing random examples, one at a time, and adjusting the learning variables using simple calculations that are usually based on the single example only. The low computational complexity (per iteration) of online algorithms is often associated with their slow convergence and low accuracy in solving the underlying optimization problems. As argued in [1, 2], the combined low complexity and low accuracy, together with other

#### Online Dictionary Learning for Sparse Coding

**Johann Mairal**  
François Bach  
INRIA,<sup>1</sup> 45 rue d’Ulm 75005 Paris, France

JULIEN.MAIRAL@INRIA.FR  
FRANCIS.BACH@INRIA.FR

**Jean Ponce**  
Ecole Normale Supérieure,<sup>2</sup> 45 rue d’Ulm 75005 Paris, France

JEAN.PONCE@ENS.FR

**Gilles Sapiro**  
University of Minnesota – Department of Electrical and Computer Engineering, 200 Union Street St., Minneapolis, USA

GILLES@TAMU.EDU

#### Abstract

Sparse coding—that is, modelling data vectors as sparse linear combinations of basis elements—is widely used in machine learning, neuroscience, signal processing, and statistics. This paper focuses on learning the basis set, also called dictionary, to adapt it to specific data, an approach that has recently proven to be very effective for signal reconstruction and classification in the audio and image processing domains. This paper proposes a new online optimization algorithm for dictionary learning, based on stochastic approximations, which scales up gracefully to large datasets with millions of training samples. A proof of convergence is presented, along with experiments with natural images demonstrating

like decompositions based on principal component analysis and its variants, these models do not impose that the basis vectors be orthogonal, allowing more flexibility to adapt the representation to the data. While learning the dictionary has proven to be critical to achieve (or improve upon) state-of-the-art results, effectively solving the corresponding optimization problem is a significant computational challenge, particularly in the context of the large-scale datasets involved in image processing tasks, that may include millions of training samples. Addressing this challenge is the topic of this paper.

Concretely, consider a signal  $\mathbf{s}$  in  $\mathbb{R}^n$ . We say that it admits a sparse approximation over a dictionary  $\mathbf{D}$  in  $\mathbb{R}^{m \times n}$ , with  $l$  columns referred to as atoms, when one can find a linear combination of a “few” atoms from  $\mathbf{D}$  that is “close” to the signal  $\mathbf{s}$ . Experiments have shown that modelling a

[Bottou & Bousquet '07]

[Xiao '09]

[Mairal et al '09]

# Outline

- 1 Back to basics: learning is optimizing
- 2 Discussion: (some) perspectives, (deep) questions, and (personal) thoughts
- 3 Highlight: fresh interest on min-max
- 4 Highlight: collaborative/federative learning

## A history of success: ML and IA

- First Generation ('90-'00): the background  
e.g., fraud detection, search engines
- Second Generation ('00-'10)  
e.g., recommendation systems
- Third Generation ('10-now): rebirth of **deep learning**  
e.g., speech recognition, computer vision, translation...

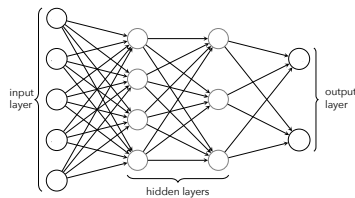
## A history of success: ML and IA

- First Generation ('90-'00): the background  
e.g., fraud detection, search engines
- Second Generation ('00-'10)  
e.g., recommendation systems
- Third Generation ('10-now): rebirth of **deep learning**  
e.g., speech recognition, computer vision, translation...
- Fourth Generation (emerging): markets  
not just one agent making decisions but multi-agents...  
towards interconnected web of data, agents, decisions

[Jordan '18] "Artificial Intelligence: The Revolution Hasn't Happened Yet"

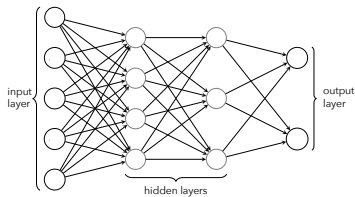
## OK but keep in mind current limits

Success of deep learning  
for image recognition

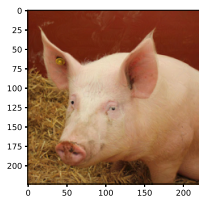


## OK but keep in mind current limits

Success of deep learning  
for image recognition



Example:

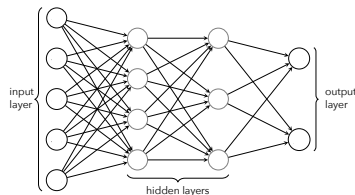


pig (99%)

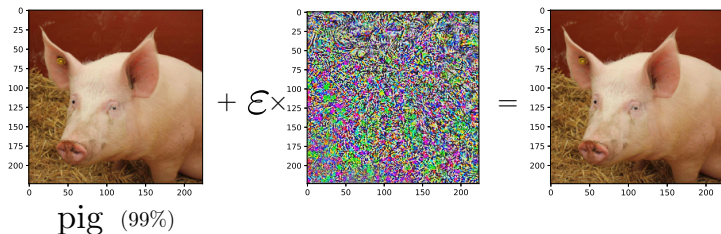


OK but keep in mind current limits

Success of deep learning  
for image recognition

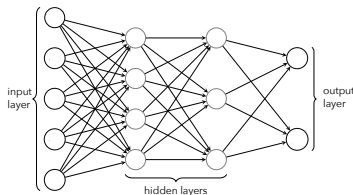


**Example:** (notebooks of NeurIPS 2018 tutorial on Adversarial Robustness)

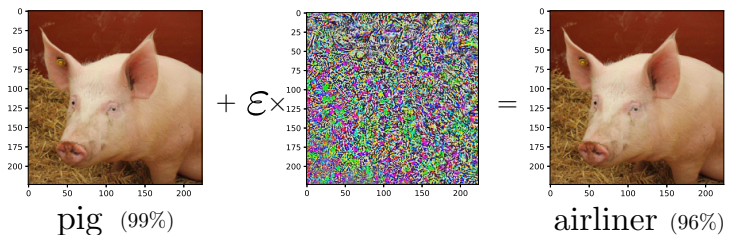


## OK but keep in mind current limits

Success of deep learning  
for image recognition



**Example:** (notebooks of NeurIPS 2018 tutorial on Adversarial Robustness)



Warning: fragile approach, and more work needed

## Beyond this: political thoughts

- ④ Deep learning needs a lot of energy



## Beyond this: political thoughts

① Deep learning needs a lot of energy



② Deep learning needs a lot of data

## Beyond this: political thoughts

- 1 Deep learning needs a lot of energy
- 2 Deep learning needs a lot of data  
(confidentiality issues)



Facebook a laissé Netflix et Spotify accéder à la messagerie privée de ses utilisateurs

Par Lucie Norbert / Modifié le 15/12/2018 à 12:01 / 146000

**Google sait où vous êtes, même si vous désactivez l'historique des positions**

Les services Google savent et stockent vos données de localisation, même si vous désactivez l'historique des positions dans vos paramètres de confidentialité.

Par Marc Dufregne avec CHATLAIN / Modifié le 14/12/2018 à 16:19 / 146000 Modifié le 14/12/2018 à 16:19

## Beyond this: political thoughts

- 1 Deep learning needs a lot of energy
- 2 Deep learning needs a lot of data  
(confidentiality issues)
- 3 Manual data treatment (in poor countries)



Facebook a laissé Netflix et Spotify accéder à la messagerie privée de ses utilisateurs

Par Lucie Norbert - Mis à jour le 15/12/2018 à 12:01 / 146000

**Google sait où vous êtes, même si vous désactivez l'historique des positions**

Les services Google savent et stockent vos données de localisation, même si vous désactivez l'historique des positions dans vos paramètres de confidentialité.

Par Marc Dufregne avec CHATIAUX - / Mis à jour le 14/12/2018 à 16:08

**amazon**  
by  
**mechanical turk**

## Beyond this: political thoughts

① Deep learning needs a lot of energy



② Deep learning needs a lot of data  
(confidentiality issues)

Facebook a laissé Netflix et Spotify accéder à la messagerie privée de ses utilisateurs

Par Lucile Norbail / Modérateur 15/12/2018 à 12:01 / 19680

**Google sait où vous êtes, même si vous désactivez l'historique des positions**

Les services Google savent et stockent vos données de localisation, même si vous désactivez l'historique des positions dans vos paramètres de confidentialité.

Par Marc Duffagni avec CNET.com / Modérateur 14 août 2018 à 16:19 / 1913 pour 1000 14 août 2018 à 16:18

③ Manual data treatment (in poor countries)



**Amazon : l'intelligence artificielle qui n'aimait pas les femmes**



Accélérer le recrutement en faisant analyser les CV par une IA : l'idée semblait prometteuse à Amazon. Mais elle s'est mise à sous-noter les femmes candidates à des postes tech.



**Machine Bias**

There's software used across the country to predict future criminals. And it's biased against blacks.

by John Angwin, MIT Review of Science and System Reviews, Published May 14, 2016

④ AI may be biased

## Beyond this: political thoughts

1 Deep learning needs a lot of energy



2 Deep learning needs a lot of data  
(confidentiality issues)

Facebook a laissé Netflix et Spotify accéder à la messagerie privée de ses utilisateurs

Par Lucile Norbail / Mis à jour le 15/12/2018 à 12:01 / 14000000

Google sait où vous êtes, même si vous désactivez l'historique des positions

Les services Google savent et stockent vos données de localisation, même si vous désactivez l'historique des positions dans vos paramètres de confidentialité.

Par Marc Duffagni avec CNRTL.com / Mis à jour le 15/12/2018 à 10:10 / 10000000

3 Manual data treatment (in poor countries)

amazon  
mechanical turk

Amazon : l'intelligence artificielle qui n'aimait pas les femmes



Accélérer le recrutement en faisant analyser les CV par une IA : l'idée semblait prometteuse à Amazon. Mais elle s'est mise à sous-noter les femmes candidates à des postes tech.



Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by John Angwin, Jeff Larson, David Rosenberg and Lauren Kirchner, ProPublica  
May 16, 2016

4 AI may be biased

5 Bad bots: fake news and massive people manipulation





## Zoom on first point: deep learning needs a lot of energy

How much **energy** spent on computational experiments of this paper ?

### CAPACITY AND TRAINABILITY IN RECURRENT NEURAL NETWORKS

Jasmine Collins; Jascha Sohl-Dickstein & David Sussillo  
Google Brain  
Google Inc.  
Mountain View, CA 94043, USA  
{j1collins, jaschasd, sussillo}@google.com

#### ABSTRACT

Two potential bottlenecks on the expressiveness of recurrent neural networks (RNNs) are their ability to store information about the task in their parameters, and to store information about the input history in their units. We show experimentally that all common RNN architectures achieve nearly the same per-task and per-unit capacity bounds with careful training, for a variety of tasks and stacking depths.

## Zoom on first point: deep learning needs a lot of energy

How much **energy** spent on computational experiments of this paper ?

Energy measure unit: (from J. Duchi, Stanford)

how many Toyota Camrys from Montpellier to Paris ? (approximately)

### CAPACITY AND TRAINABILITY IN RECURRENT NEURAL NETWORKS

Jasmine Collins,<sup>\*</sup> Jascha Sohl-Dickstein & David Sussillo  
Google Brain  
Google Inc.  
Mountain View, CA 94043, USA  
{jcollins, jaschasd, sussillo}@google.com

#### ABSTRACT

Two potential bottlenecks on the expressiveness of recurrent neural networks (RNNs) are their ability to store information about the task in their parameters, and to store information about the input history in their units. We show experimentally that all common RNN architectures achieve nearly the same per-task and per-unit capacity bounds with careful training, for a variety of tasks and stacking depths.



1?

## Zoom on first point: deep learning needs a lot of energy

How much **energy** spent on computational experiments of this paper ?

Energy measure unit: (from J. Duchi, Stanford)

how many Toyota Camrys from Montpellier to Paris ? (approximately)

### CAPACITY AND TRAINABILITY IN RECURRENT NEURAL NETWORKS

Jasmine Collins, Jascha Sohl-Dickstein & David Sussillo  
Google Brain  
Mountain View, CA 94043, USA  
{jcollins, jaschasd, sussillo}@google.com

#### ABSTRACT

Two potential bottlenecks on the expressiveness of recurrent neural networks (RNNs) are their ability to store information about the task in their parameters, and to store information about the input history in their units. We show experimentally that all common RNN architectures achieve nearly the same per-task and per-unit capacity bounds with careful training, for a variety of tasks and stacking depths.



10?

## Zoom on first point: deep learning needs a lot of energy

How much **energy** spent on computational experiments of this paper ?

Energy measure unit: (from J. Duchi, Stanford)

how many Toyota Camrys from Montpellier to Paris ? (approximately)

### CAPACITY AND TRAINABILITY IN RECURRENT NEURAL NETWORKS

Jasmine Collins; Jascha Sohl-Dickstein & David Sussillo  
Google Brain  
Google Inc.  
Mountain View, CA 94043, USA  
{jcollins, jaschasd, sussillo}@google.com

#### ABSTRACT

Two potential bottlenecks on the expressiveness of recurrent neural networks (RNNs) are their ability to store information about the task in their parameters, and to store information about the input history in their units. We show experimentally that all common RNN architectures achieve nearly the same per-task and per-unit capacity bounds with careful training, for a variety of tasks and stacking depths.



100?

## Zoom on first point: deep learning needs a lot of energy

How much **energy** spent on computational experiments of this paper ?

Energy measure unit: (from J. Duchi, Stanford)

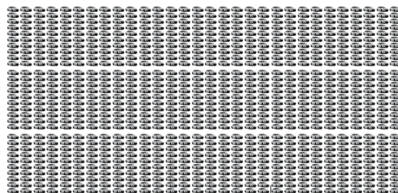
how many Toyota Camrys from Montpellier to Paris ? (approximately)

### CAPACITY AND TRAINABILITY IN RECURRENT NEURAL NETWORKS

Jasmine Collins, Jascha Sohl-Dickstein & David Sussillo  
Google Brain  
Mountain View, CA 94043, USA  
{jcollins, jaschasd, sussillo}@google.com

#### ABSTRACT

Two potential bottlenecks on the expressiveness of recurrent neural networks (RNNs) are their ability to store information about the task in their parameters, and to store information about the input history in their units. We show experimentally that all common RNN architectures achieve nearly the same per-task and per-unit capacity bounds with careful training, for a variety of tasks and stacking depths.



1000?

## Zoom on first point: deep learning needs a lot of energy

How much **energy** spent on computational experiments of this paper ?

Energy measure unit: (from J. Duchi, Stanford)

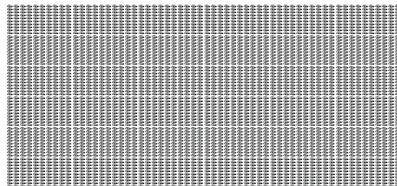
how many Toyota Camrys from Montpellier to Paris ? (approximately)

### CAPACITY AND TRAINABILITY IN RECURRENT NEURAL NETWORKS

Jasmine Collins\*, Jascha Sohl-Dickstein & David Sussillo  
Google Brain  
Google Inc.  
Mountain View, CA 94043, USA  
{jicollins, jaschasd, sussillo}@google.com

#### ABSTRACT

Two potential bottlenecks on the expressiveness of recurrent neural networks (RNNs) are their ability to store information about the task in their parameters, and to store information about the input history in their units. We show experimentally that all common RNN architectures achieve nearly the same per-task and per-unit capacity bounds with careful training, for a variety of tasks and stacking depths.



4200

Deep learning has a terrible carbon footprint...

→ work needed towards energy-efficient learning models

→ for us, in particular: efficient optimization algorithms

(fast convergence, automatic tuning of parameters,...)

# Outline

- 1 Back to basics: learning is optimizing
- 2 Discussion: (some) perspectives, (deep) questions, and (personal) thoughts
- 3 Highlight: fresh interest on min-max
- 4 Highlight: collaborative/federative learning

## Min-Max & friends

In OR/Optimization/Games,

we are used to deal with min-max, (Nash) equilibrium, or saddle-points...

$$\min_{x \in X} \max_{y \in Y} F(x, y) \quad \text{or} \quad F(x^*, y) \leq F(x^*, y^*) \leq F(x, y^*)$$

Examples: nonconvex problems (with  $X$  discrete)

- Lagrangian duality or relaxation

e.g. [Lemaréchal '01] “the omnipresence of Lagrange”

$$\min_u \max_{x \in X} L(x, u) = p^\top x - u^\top c(x)$$

- Robust optimization

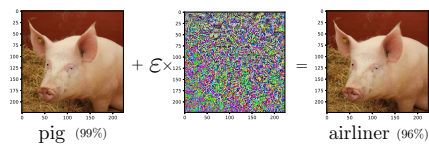
e.g. the work of our conference chair [Poss '18], [Poss et al '16]

$$\min_{x \in X} \max_{\xi \in \Delta_x} f(x, \xi)$$



# Examples in AI: non-convex too

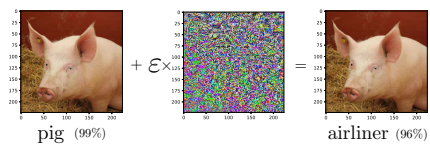
## Adversarially robust models [Kolter & Madry '18]



$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \max_{\delta \in \Delta} \ell(y_i, h(x, a_i + \delta))$$

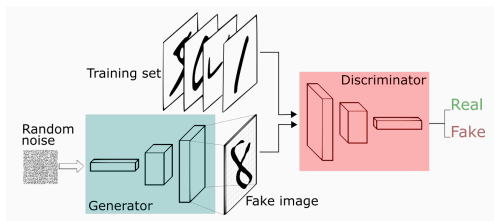
## Examples in AI: non-convex too

### Adversarially robust models [Kolter & Madry '18]



$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \max_{\delta \in \Delta} \ell(y_i, h(x, a_i + \delta))$$

### Generate data with GANs [Goodfellow et al '14]



$$\min_{x_G} \max_{x_D} \mathbb{E}_{a \sim \text{data}} \left[ \log h_D(x_D, a) \right] + \mathbb{E}_{z \sim \text{noise}} \left[ \log \left( 1 - h_D(x_D, h_G(x_G, z)) \right) \right]$$

## GANs: sucesses and failures

Question: who is real, who isn't ?



## GANs: sucesses and failures

Question: who is real, who isn't ?

Answer: **both** are fake !

[<https://thispersondoesnotexist.com>]



## GANs: sucesses and failures

Question: who is real, who isn't ?

Answer: **both** are fake !

[<https://thispersondoesnotexist.com>]



But the story far from being over...

## GANs: successes and failures

Question: who is real, who isn't ?

Answer: **both** are fake !

[<https://thispersondoesnotexist.com>]



But the story far from being over...

Coupling of two neural networks  
gives rise to strange behaviors

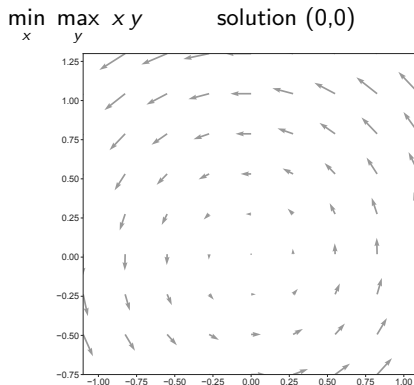
Even when solved with  
state-of-the-art stochastic  
gradient (extra-gradient variants)



## Example of strange phenomena... and a simple fix

Non-convergent phenomena are observed even in very simple problems

Example:

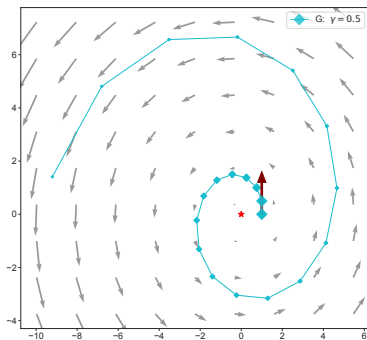


## Example of strange phenomena... and a simple fix

Non-convergent phenomena are observed even in very simple problems

Example:

min  $x$    max  $y$    solution (0,0)



- Gradient algorithm diverges...

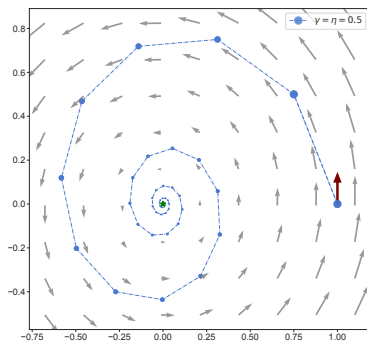


## Example of strange phenomena... and a simple fix

Non-convergent phenomena are observed even in very simple problems

Example:

min max  $x$   $y$  solution (0,0)



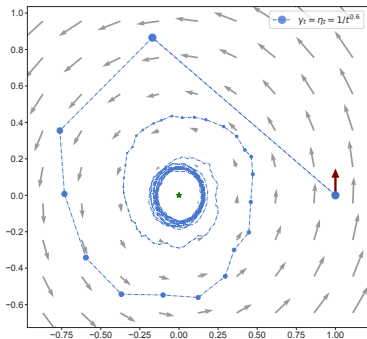
- Gradient algorithm diverges...
- Extra-gradient algorithm converges [Korpelevich '76] ( → GANs)

## Example of strange phenomena... and a simple fix

Non-convergent phenomena are observed even in very simple problems

Example:

min  $x$    max  $y$    solution (0,0)

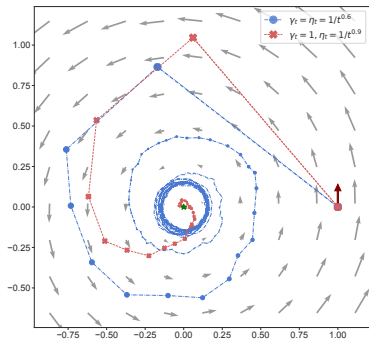


- Gradient algorithm diverges...
- Extra-gradient algorithm converges [Korpelevich '76] ( → GANs)
- Stochastic extra-gradient never converges...

## Example of strange phenomena... and a simple fix

Non-convergent phenomena are observed even in very simple problems

Example:  $\min_x \max_y$  solution (0,0)



- Gradient algorithm diverges...
- Extra-gradient algorithm converges [Korpelevich '76] ( $\rightarrow$  GANs)
- Stochastic extra-gradient never converges...
- A remedy: use double stepsize [Hsieh, Iutzeler, M., Mertikopoulos '20]

# Outline

- 1 Back to basics: learning is optimizing
- 2 Discussion: (some) perspectives, (deep) questions, and (personal) thoughts
- 3 Highlight: fresh interest on min-max
- 4 Highlight: collaborative/federative learning

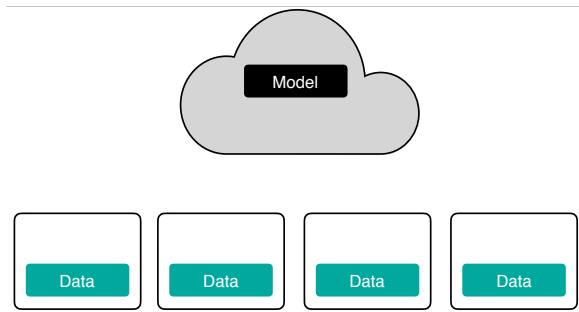
# Collaborative learning

Set-up:



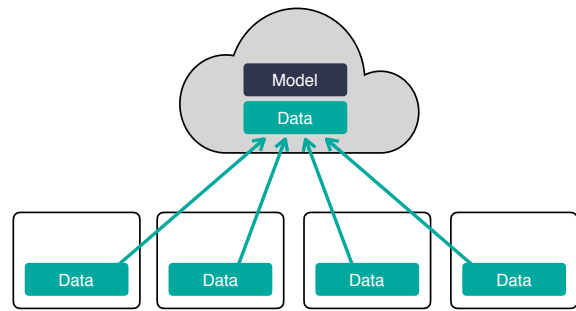
# Collaborative learning

Set-up: (standard) centralized learning



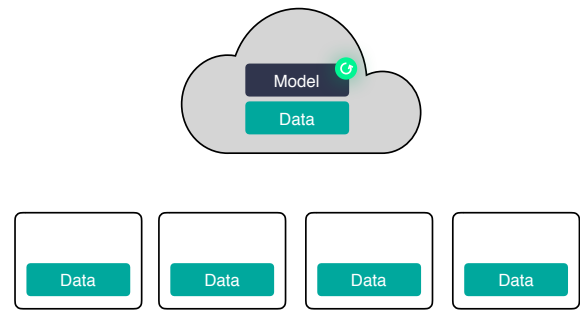
# Collaborative learning

Set-up: (standard) centralized learning



# Collaborative learning

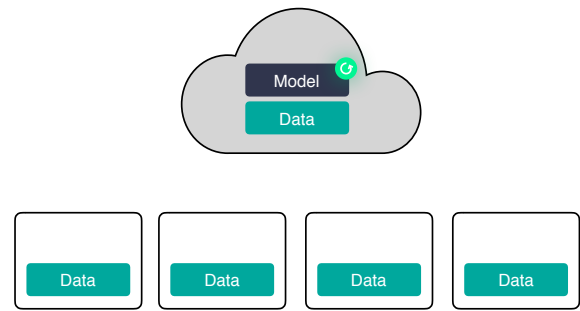
Set-up: (standard) centralized learning





# Collaborative learning

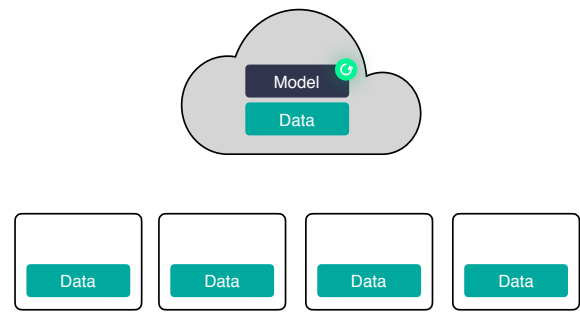
Set-up: (standard) centralized learning



- needs of lot of storage 😞... but efficient 😊

# Collaborative learning

Set-up: (standard) centralized learning

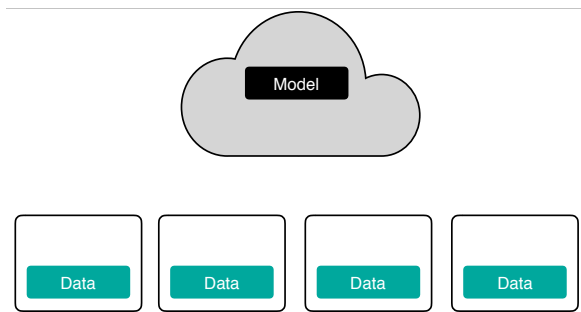


- needs of lot of storage 😞... but efficient 😊
- is highly privacy invasive (e.g. phones) 😞
- jeopardy on confidentiality/strategy (e.g. hospitals/companies) 😞

# Move the model, not the data

## Different set-up: Federative Learning

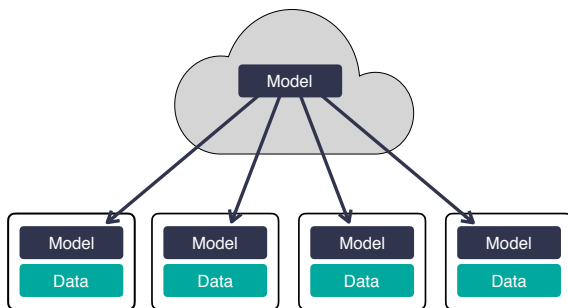
Decoupling the ability to learn a global model from moving local data



# Move the model, not the data

## Different set-up: Federative Learning

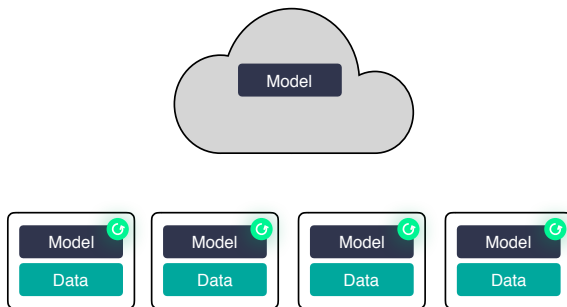
Decoupling the ability to learn a global model from moving local data



## Move the model, not the data

### Different set-up: Federative Learning

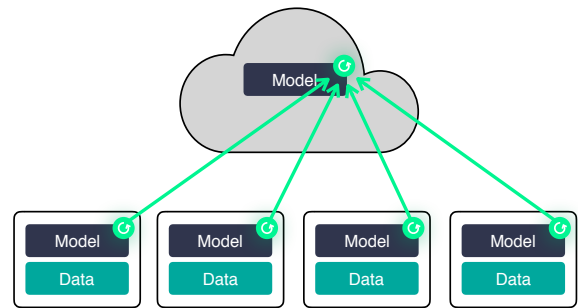
Decoupling the ability to learn a global model from moving local data



## Move the model, not the data

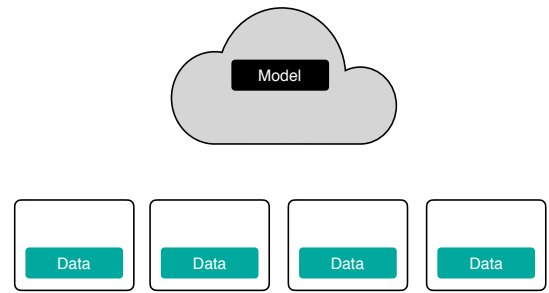
### Different set-up: Federative Learning

Decoupling the ability to learn a global model from moving local data

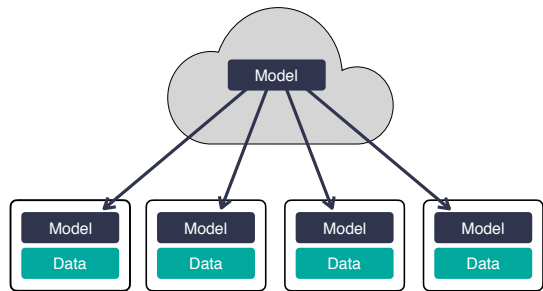


New (optimization) algorithm able to deal with **heterogenous** data/systems

## Zoom on a research topic: harnessing communications

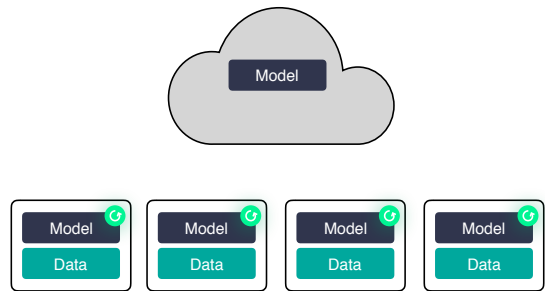


## Zoom on a research topic: harnessing communications

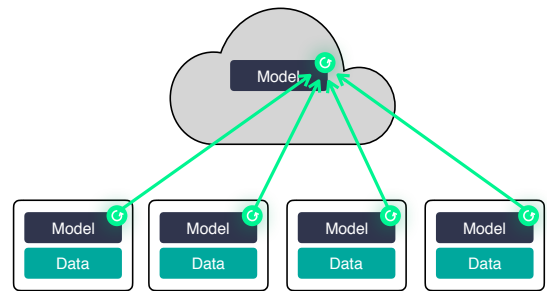




## Zoom on a research topic: harnessing communications

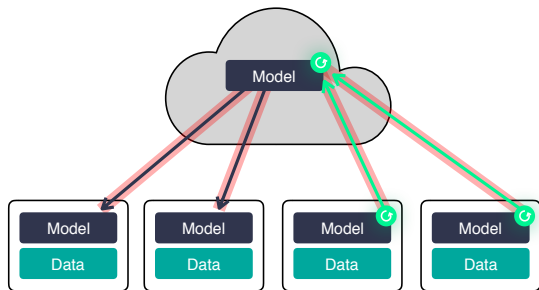


## Zoom on a research topic: harnessing communications



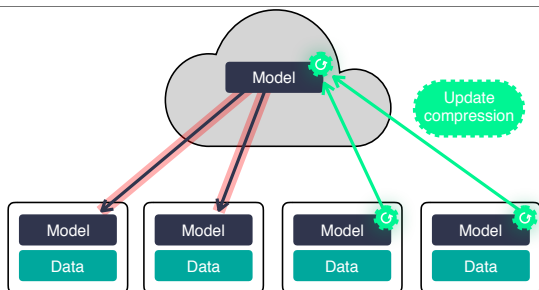
## Zoom on a research topic: harnessing communications

Communication is the bottleneck 😞



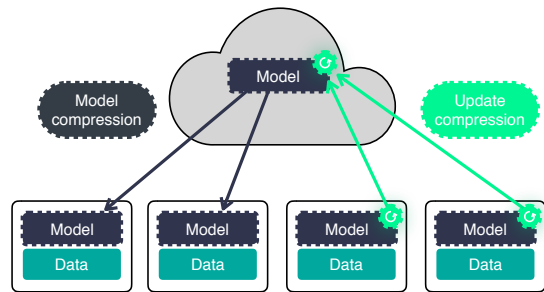
## Zoom on a research topic: harnessing communications

Communication is the bottleneck 😞



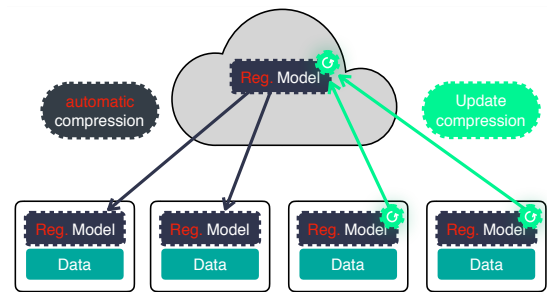
## Zoom on a research topic: harnessing communications

Communication is the bottleneck 😞



## Zoom on a research topic: harnessing communications

Communication is the bottleneck 😞

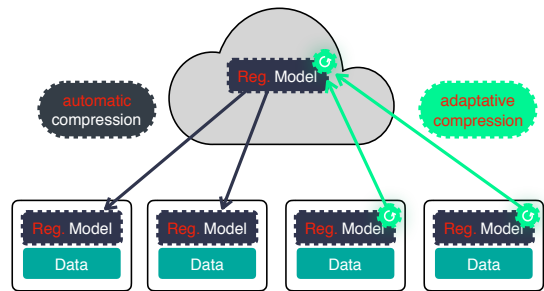


Our contribution: compression by nonsmooth regularization

- Observation: nonsmooth regularization gives automatic model compression  
E.g. for  $r = \|\cdot\|_1$ , model becomes sparse... just communicate nonzeros!

## Zoom on a research topic: harnessing communications

Communication is the bottleneck 😞



### Our contribution: compression by nonsmooth regularization

- Observation: nonsmooth regularization gives automatic model compression  
E.g. for  $r = \|\cdot\|_1$ , model becomes sparse... just communicate nonzeros!
- [Grishchenko, Iutzeler, M. '19] uses it for update comp.  
E.g. for  $r = \|\cdot\|_1$ , select current support + random entries

# Illustration of compression by a nonsmooth regularizer

On an instance of TV-regularized logistic regression (a1a dataset on 10 machines)

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{j=1}^n \log(1 + \exp(-y_j a_j^\top x)) + \lambda \text{TV}(x)$$

$$\text{TV}(x) = \sum_{i=1}^{n-1} |x_{i+1} - x_i|$$

Total Variation

- Comparison of
- Usual algorithm (black)
  - Our variant with compression (red)



# Illustration of compression by a nonsmooth regularizer

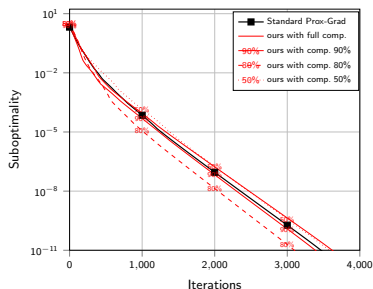
On an instance of TV-regularized logistic regression (a1a dataset on 10 machines)

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{j=1}^n \log(1 + \exp(-y_j a_j^\top x)) + \lambda \text{TV}(x)$$

$$\text{TV}(x) = \sum_{i=1}^{n-1} |x_{i+1} - x_i|$$

Total Variation

- Comparison of
- Usual algorithm (black)
  - Our variant with compression (red)



# Illustration of compression by a nonsmooth regularizer

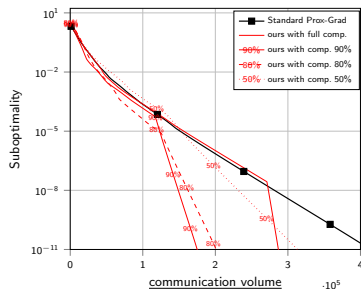
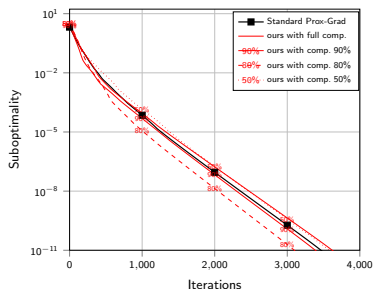
On an instance of TV-regularized logistic regression (a1a dataset on 10 machines)

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{j=1}^n \log(1 + \exp(-y_j a_j^\top x)) + \lambda \text{TV}(x)$$

$$\text{TV}(x) = \sum_{i=1}^{n-1} |x_{i+1} - x_i|$$

Total Variation

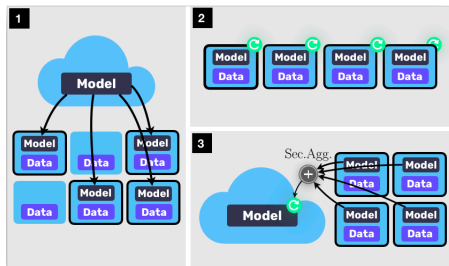
- Comparison of
- Usual algorithm (black)
  - Our variant with compression (red)



Acceleration... with respect to communication volume !

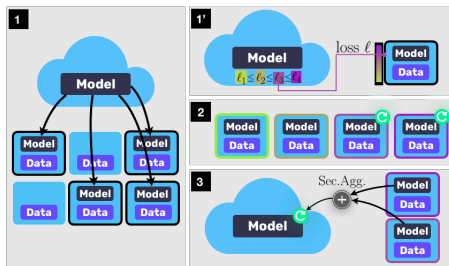
## Zoom on a second topic: improving worst-case FedAvg

Federated Learning by Google = FedAvg



## Zoom on a second topic: improving worst-case FedAvg

Federated Learning by Google = FedAvg vs Robust FedAvg

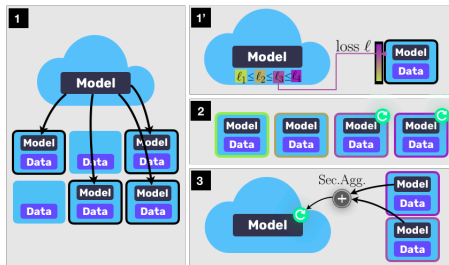


Our contribution: improve worst-case performance over users

- by adaptive filtering [Laguel, M., Harchaoui '19]

## Zoom on a second topic: improving worst-case FedAvg

Federated Learning by Google = **FedAvg** vs **Robust FedAvg**



Our contribution: improve worst-case performance over users

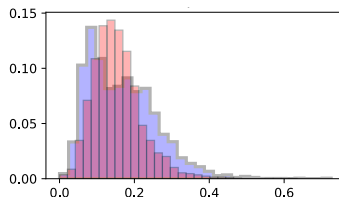
- by adaptive filtering [Laguel, M., Harchaoui '19]

Example: classification task

**Standard** vs **Ours**

(using ConvNet on EMNIST dataset)

Histogram over users  
of test misclassification error



# Conclusion

## Take-home message

- AI is a rich, active, visible field of research and developments
- Next step: Operational AI (?) (not much discussed today... except in the title!)
- There are many problems involving uncertainty, decision-making, robustness and scale... far from being solved  
(not to mention economic, social and legal issues...)

## Two theoretical questions

- understanding convergence towards good equilibrium
- models and algorithms in federated learning  
compression by nonsmooth regularization, improvement of worst-case

# Conclusion

## Take-home message

- AI is a rich, active, visible field of research and developments
- Next step: Operational AI (?) (not much discussed today... except in the title!)
- There are many problems involving uncertainty, decision-making, robustness and scale... far from being solved  
(not to mention economic, social and legal issues...)

## Two theoretical questions

- understanding convergence towards good equilibrium
- models and algorithms in federated learning  
compression by nonsmooth regularization, improvement of worst-case

thanks !!