

Using Integer Programming to Train Neural Networks

Daniel Bienstock, Gonzalo Muñoz, Sebastian Pokutta

Columbia University, Polytechnique Montréal, Georgia Tech

Bordeaux July 2018

Empirical Risk Minimization problem

Empirical Risk Minimization problem

Given:

- D data points $(\mathbf{x}_i, \mathbf{y}_i)$, $i = 1, \dots, D$
- $\mathbf{x}_i \in \mathbb{R}^n$, $\mathbf{y}_i \in \mathbb{R}^m$ all i
- A “loss” function $\ell : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ (not necessarily convex)

Compute $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ to solve

$$\min_{\mathbf{f}} \frac{1}{D} \sum_{i=1}^D \ell(\mathbf{f}(\mathbf{x}^i), \mathbf{y}^i) \quad (+ \text{ optional regularizer } \Phi(\mathbf{f}))$$

$$\mathbf{f} \in \mathbf{F} \quad (\text{some class})$$

Approximate optimization of well-behaved functions

Approximate optimization of well-behaved functions

Prototype problem:

$$\begin{aligned} & \min c^T x \\ \text{s.t. } & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & x \in [0, 1]^n \end{aligned}$$

Approximate optimization of well-behaved functions

Prototype problem:

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & x \in [0, 1]^n \end{aligned}$$

Each f_i is “well-behaved”: Lipschitz constant \mathcal{L}_i

Approximate optimization of well-behaved functions

Prototype problem:

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & x \in [0, 1]^n \end{aligned}$$

Each f_i is “well-behaved”: Lipschitz constant \mathcal{L}_i

Toolset:

- **Intersection graph**

Approximate optimization of well-behaved functions

Prototype problem:

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & x \in [0, 1]^n \end{aligned}$$

Each f_i is “well-behaved”: Lipschitz constant \mathcal{L}_i

Toolset:

- **Intersection graph**

A vertex for each variable and an edge whenever two variables appear in the same f_i

- **Tree-width**

Approximate optimization of well-behaved functions

Prototype problem:

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & x \in [0, 1]^n \end{aligned}$$

Each f_i is “well-behaved”: Lipschitz constant \mathcal{L}_i

Toolset:

- **Intersection graph**

A vertex for each variable and an edge whenever two variables appear in the same f_i

- **Tree-width** Min clique number (minus one) over all chordal supergraphs of G

Prototype problem:

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & x \in [0, 1]^n \end{aligned}$$

An extension of work in B. and Muñoz 2015, SIOPT 2018.

Suppose:

the intersection graph has tree-width ω and f_i has Lipschitz constant $\mathcal{L}_i \leq \mathcal{L}$.

Prototype problem:

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & x \in [0, 1]^n \end{aligned}$$

An extension of work in B. and Muñoz 2015, SIOPT 2018.

Suppose:

the intersection graph has tree-width ω and f_i has Lipschitz constant $\mathcal{L}_i \leq \mathcal{L}$.

Then, for every $0 < \epsilon < 1$ there is an **LP** relaxation with

Prototype problem:

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & x \in [0, 1]^n \end{aligned}$$

An extension of work in B. and Muñoz 2015, SIOPT 2018.

Suppose:

the intersection graph has tree-width ω and f_i has Lipschitz constant $\mathcal{L}_i \leq \mathcal{L}$.

Then, for every $0 < \epsilon < 1$ there is an **LP** relaxation with

$O((\mathcal{L}/\epsilon)^{\omega+1} n \log(\mathcal{L}/\epsilon))$ variables and constraints

Prototype problem:

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & x \in [0, 1]^n \end{aligned}$$

An extension of work in B. and Muñoz 2015, SIOPT 2018.

Suppose:

the intersection graph has tree-width ω and f_i has Lipschitz constant $\mathcal{L}_i \leq \mathcal{L}$.

Then, for every $0 < \epsilon < 1$ there is an **LP** relaxation with

$O((\mathcal{L}/\epsilon)^{\omega+1} n \log(\mathcal{L}/\epsilon))$ variables and constraints

Optimality and feasibility errors $O(\epsilon)$

Application to ERM problem

$$\min_{f \in \mathcal{F}} \frac{1}{D} \sum_{i=1}^D \ell(f(x^i), y^i)$$

Application to ERM problem

$$\min_{f \in \mathcal{F}} \frac{1}{D} \sum_{i=1}^D \ell(f(x^i), y^i)$$

Linearize objective using epigraph trick

$$\begin{aligned} \min_{f \in \mathcal{F}} \frac{1}{D} \sum_{i=1}^D L_i \\ L_i \geq \ell(f(x^i), y^i) \quad 1 \leq i \leq D \end{aligned}$$

Function parameterization

$$\min_{f \in \mathcal{F}} \frac{1}{D} \sum_{i=1}^D L_i$$
$$L_i \geq \ell(f(x^i), y^i) \quad 1 \leq i \leq D$$

Examples:

- **Neural Networks with k layers.**

$f(x) = T_{k+1} \circ \sigma \circ T_k \circ \sigma \dots \circ \sigma \circ T_1(x)$, each T_j affine, parameterized by its coefficients.

- **Linear Regression.** $f(x) = Ax + b$ with ℓ_2 -loss.
- **Binary Classification.** Varying f architectures and cross-entropy loss: $\ell(p, y) = -y \log(p) - (1 - y) \log(1 - p)$

Function parameterization

$$\min_{f \in F} \frac{1}{D} \sum_{i=1}^D L_i$$
$$L_i \geq \ell(f(x^i), y^i) \quad 1 \leq i \leq D$$

Examples:

- **Neural Networks with k layers.**
 $f(x) = T_{k+1} \circ \sigma \circ T_k \circ \sigma \dots \circ \sigma \circ T_1(x)$, each T_j affine, parameterized by its coefficients.
- **Linear Regression.** $f(x) = Ax + b$ with ℓ_2 -loss.
- **Binary Classification.** Varying f architectures and cross-entropy loss: $\ell(p, y) = -y \log(p) - (1 - y) \log(1 - p)$

We assume $F = \{f(x, \theta) \text{ for } \theta \in \Theta \subseteq [-1, 1]^N\}$

We now apply the LP approximation result to:

$$\min_{\theta \in \Theta} \frac{1}{D} \sum_{i=1}^D L_i$$
$$L_i \geq \ell(f(x^i, \theta), y^i) \quad 1 \leq i \leq D$$

Let θ^* be an optimal solution. For every $0 < \epsilon < 1$ there is an **LP** of size

$$O\left(\left(\mathcal{L}/\epsilon\right)^{2N+2} (N + D) \log(\mathcal{L}/\epsilon)\right)$$

We now apply the LP approximation result to:

$$\begin{aligned} \min_{\theta \in \Theta} \frac{1}{D} \sum_{i=1}^D L_i \\ L_i \geq \ell(f(x^i, \theta), y^i) \quad 1 \leq i \leq D \end{aligned}$$

Let θ^* be an optimal solution. For every $0 < \epsilon < 1$ there is an **LP** of size

$$O\left(\left(\mathcal{L}/\epsilon\right)^{2N+2} (N + D) \log(\mathcal{L}/\epsilon)\right)$$

such that its optimal solution $\hat{\theta}$ satisfies:

$$\frac{1}{D} \sum_{i=1}^D \ell(f(x^i, \hat{\theta}), y^i) \leq \frac{1}{D} \sum_{i=1}^D \ell(f(x^i, \theta^*), y^i) + \epsilon$$

We now apply the LP approximation result to:

$$\begin{aligned} \min_{\theta \in \Theta} \frac{1}{D} \sum_{i=1}^D L_i \\ L_i \geq \ell(f(x^i, \theta), y^i) \quad 1 \leq i \leq D \end{aligned}$$

Let θ^* be an optimal solution. For every $0 < \epsilon < 1$ there is an **LP** of size

$$O((\mathcal{L}/\epsilon)^{2N+2} (N + D) \log(\mathcal{L}/\epsilon))$$

such that its optimal solution $\hat{\theta}$ satisfies:

$$\frac{1}{D} \sum_{i=1}^D \ell(f(x^i, \hat{\theta}), y^i) \leq \frac{1}{D} \sum_{i=1}^D \ell(f(x^i, \theta^*), y^i) + \epsilon$$

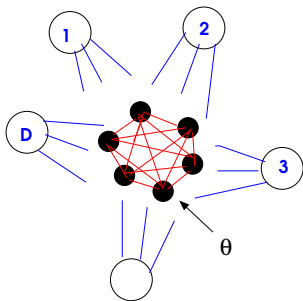
\mathcal{L} is an upper bound on the Lipschitz constant of $g^i(\theta) \doteq \ell(f(x^i, \theta), y^i)$.

Proof sketch:

$$\min_{\theta \in \Theta} \frac{1}{D} \sum_{i=1}^D L_i, \quad \text{s.t.} \quad L_i \geq \ell(f(x^i, \theta), y^i) \quad 1 \leq i \leq D$$

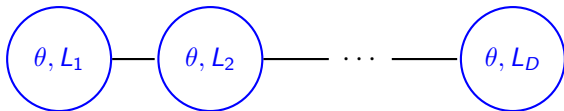
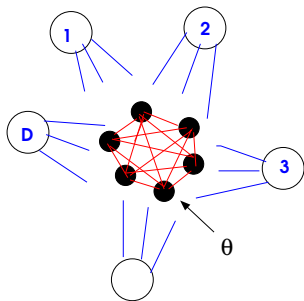
Proof sketch:

$$\min_{\theta \in \Theta} \frac{1}{D} \sum_{i=1}^D L_i, \quad \text{s.t.} \quad L_i \geq \ell(f(x^i, \theta), y^i) \quad 1 \leq i \leq D$$



Proof sketch:

$$\min_{\theta \in \Theta} \frac{1}{D} \sum_{i=1}^D L_i, \quad \text{s.t.} \quad L_i \geq \ell(f(x^i, \theta), y^i) \quad 1 \leq i \leq D$$



Application: training of Deep Neural Networks with ReLUs

Application: training of Deep Neural Networks with ReLUs

As per Arora Basu Mianjy Mukherjee ICLR '18

The setup:

- D data points (\mathbf{x}_i, y_i) , $1 \leq i \leq D$, $\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$

Application: training of Deep Neural Networks with ReLUs

As per Arora Basu Mianjy Mukherjee ICLR '18

The setup:

- D data points (\mathbf{x}_i, y_i) , $1 \leq i \leq D$, $\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$
- **Task: compute** a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ to minimize

$$\frac{1}{D} \sum_{i=1} (y_i - f(\mathbf{x}_i))^2$$

Application: training of Deep Neural Networks with ReLUs

As per Arora Basu Mianjy Mukherjee ICLR '18

The setup:

- D data points (\mathbf{x}_i, y_i) , $1 \leq i \leq D$, $\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$
- **Task: compute** a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ to minimize

$$\frac{1}{D} \sum_{i=1}^D (y_i - f(\mathbf{x}_i))^2$$

- $f = T_{k+1} \circ \sigma \circ T_k \circ \sigma \dots \circ \sigma \circ T_1$ (“ \circ ” = composition)

Application: training of Deep Neural Networks with ReLUs

As per Arora Basu Mianjy Mukherjee ICLR '18

The setup:

- D data points (\mathbf{x}_i, y_i) , $1 \leq i \leq D$, $\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$
- **Task: compute** a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ to minimize

$$\frac{1}{D} \sum_{i=1} (y_i - f(\mathbf{x}_i))^2$$

- $f = T_{k+1} \circ \sigma \circ T_k \circ \sigma \dots \circ \sigma \circ T_1$ (“ \circ ” = composition)
- $\sigma(t) = \max\{0, t\}$

Application: training of Deep Neural Networks with ReLUs

As per Arora Basu Mianjy Mukherjee ICLR '18

The setup:

- D data points (\mathbf{x}_i, y_i) , $1 \leq i \leq D$, $\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$
- **Task: compute** a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ to minimize

$$\frac{1}{D} \sum_{i=1} (y_i - f(\mathbf{x}_i))^2$$

- $f = T_{k+1} \circ \sigma \circ T_k \circ \sigma \dots \circ \sigma \circ T_1$ (“ \circ ” = composition)
- $\sigma(t) = \max\{0, t\}$
- Each T_h affine: $T_h(\mathbf{y}) = \mathbf{A}_h \mathbf{y} + \mathbf{b}_h$.

Application: training of Deep Neural Networks with ReLUs

As per Arora Basu Mianjy Mukherjee ICLR '18

The setup:

- D data points (\mathbf{x}_i, y_i) , $1 \leq i \leq D$, $\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$
- **Task:** compute a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ to minimize

$$\frac{1}{D} \sum_{i=1} (y_i - f(\mathbf{x}_i))^2$$

- $f = T_{k+1} \circ \sigma \circ T_k \circ \sigma \dots \circ \sigma \circ T_1$ (“ \circ ” = composition)
- $\sigma(t) = \max\{0, t\}$
- Each T_h affine: $T_h(\mathbf{y}) = \mathbf{A}_h \mathbf{y} + \mathbf{b}_h$,
- For some \mathbf{w} , \mathbf{A}_1 is $n \times \mathbf{w}$, \mathbf{A}_{k+1} is $\mathbf{w} \times \mathbf{1}$, \mathbf{A}_h is $\mathbf{w} \times \mathbf{w}$ otherwise.

Application: training of Deep Neural Networks with ReLUs

As per Arora Basu Mianjy Mukherjee ICLR '18

The setup:

- D data points (\mathbf{x}_i, y_i) , $1 \leq i \leq D$, $\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$
- **Task:** compute a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ to minimize

$$\frac{1}{D} \sum_{i=1} (y_i - f(\mathbf{x}_i))^2$$

- $f = T_{k+1} \circ \sigma \circ T_k \circ \sigma \dots \circ \sigma \circ T_1$ (“ \circ ” = composition)
- $\sigma(t) = \max\{0, t\}$
- Each T_h affine: $T_h(\mathbf{y}) = \mathbf{A}_h \mathbf{y} + \mathbf{b}_h$,
- For some \mathbf{w} , \mathbf{A}_1 is $n \times \mathbf{w}$, \mathbf{A}_{k+1} is $\mathbf{w} \times \mathbf{1}$, \mathbf{A}_h is $\mathbf{w} \times \mathbf{w}$ otherwise. Similarly with the \mathbf{b}_h .

- D data points (x_i, y_i) , $1 \leq i \leq D$, $x_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$
- **Task:** compute a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ to minimize

$$\frac{1}{D} \sum_{i=1}^D (y_i - f(x_i))^2$$

- $f = T_{k+1} \circ \sigma \circ T_k \circ \sigma \dots \circ \sigma \circ T_1$ (“ \circ ” = composition)
- $\sigma(t) = \max\{0, t\}$
- Each T_h affine: $T_h(y) = A_h y + b_h$,
- For some w , A_1 is $n \times w$, A_{k+1} is $w \times 1$, A_h is $w \times w$ otherwise. Similarly with the b_h .

- D data points (x_i, y_i) , $1 \leq i \leq D$, $x_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$
- **Task: compute** a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ to minimize

$$\frac{1}{D} \sum_{i=1}^D (y_i - f(x_i))^2$$

- $f = T_{k+1} \circ \sigma \circ T_k \circ \sigma \dots \circ \sigma \circ T_1$ (“ \circ ” = composition)
- $\sigma(t) = \max\{0, t\}$
- Each T_h affine: $T_h(y) = A_h y + b_h$,
- For some w , A_1 is $n \times w$, A_{k+1} is $w \times 1$, A_h is $w \times w$ otherwise. Similarly with the b_h .

Theorem (Arora et al 2018).

If $k = 1$ (one “hidden layer”) there is an exact algorithm of complexity

- D data points (x_i, y_i) , $1 \leq i \leq D$, $x_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$
- **Task:** compute a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ to minimize

$$\frac{1}{D} \sum_{i=1}^D (y_i - f(x_i))^2$$

- $f = T_{k+1} \circ \sigma \circ T_k \circ \sigma \dots \circ \sigma \circ T_1$ (“ \circ ” = composition)
- $\sigma(t) = \max\{0, t\}$
- Each T_h affine: $T_h(y) = A_h y + b_h$,
- For some w , A_1 is $n \times w$, A_{k+1} is $w \times 1$, A_h is $w \times w$ otherwise. Similarly with the b_h .

Theorem (Arora et al 2018).

If $k = 1$ (one “hidden layer”) there is an exact algorithm of complexity

$$O(2^w D^{nw} \text{poly}(D, n, w))$$

Polynomial in the size of the data set, for fixed n, w

Our result: if the entries of A_i, b_i are required to be in $[-1, 1]$, for any k, n, w, ϵ there is an LP of size

$$O\left(\left(\frac{w}{\epsilon}\right)^{\text{poly}(n,k,w)} (\text{poly}(n, k, w) + D) \log(w/\epsilon)\right)$$

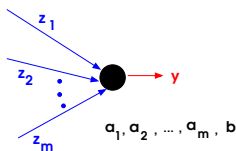
Our result: if the entries of A_i, b_i are required to be in $[-1, 1]$, for any k, n, w, ϵ there is an LP of size

$$O\left(\left(w/\epsilon\right)^{\text{poly}(n,k,w)} \left(\text{poly}(n, k, w) + D\right) \log(w/\epsilon)\right)$$

- $\text{poly}(n, k, w) =$ quadratic in w , in k , linear in n
- Treewidth **independent** of D
- Number of variables **linear** in D

The Arora-Blum setup (Binarized Neural Networks)

Activation units:

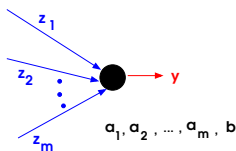


With $z \in \{0, 1\}^m$,

$$y = \begin{cases} 1, & \text{if } a^T z > b \\ 0, & \text{otherwise.} \end{cases}$$

The Arora-Blum setup (Binarized Neural Networks)

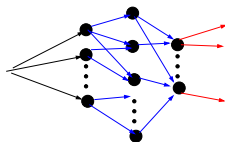
Activation units:



With $z \in \{0, 1\}^m$,

$$y = \begin{cases} 1, & \text{if } a^T z > b \\ 0, & \text{otherwise.} \end{cases}$$

Network with n binary inputs, m binary outputs, k layers



Binarized Neural Networks, 2

Training data: Set of D pairs (x^i, y^i) , $1 \leq i \leq D$
 $x^i \in \{0, 1\}^n$, $y^i \in \{0, 1\}^m$

Binarized Neural Networks, 2

Training data: Set of D pairs (x^i, y^i) , $1 \leq i \leq D$
 $x^i \in \{0, 1\}^n$, $y^i \in \{0, 1\}^m$

Problem: Compute the activation function at each node to

$$\min \frac{1}{D} \sum_{i=1}^D \ell(f(x^i), y^i)$$

(f = network function)

- When $\ell \in$ (absolute value, 2-norm squared) NP-hard if $k = 3$ and $D = n$, $m = 1$
- But we are interested in the case D very large compared to n
- And also other loss functions, e.g. smooth convex

Our result on Binarized Neural Networks

Training data: Set of D pairs (x^i, y^i) , $1 \leq i \leq D$
 $x^i \in \{0, 1\}^n$, $y^i \in \{0, 1\}^m$

Our result on Binarized Neural Networks

Training data: Set of D pairs (x^i, y^i) , $1 \leq i \leq D$
 $x^i \in \{0, 1\}^n$, $y^i \in \{0, 1\}^m$

ERM problem: Compute the activation function at each node to

$$\min \frac{1}{D} \sum_{i=1}^D \ell(f(x^i), y^i)$$

(f = network function)

Theorem. When $\ell \in$ (absolute value, 2-norm squared) there is an LP of encoding length

$$O(2^{\text{poly}(k,n,m)}(\text{poly}(k, n, m) + D))$$

that solves the ERM problem exactly in absolute case, and within $O(\epsilon)$ additive error in the 2-norm case.

Extensions