



# Strong mixed-integer programming formulations for trained neural networks

### Joey Huchette<sup>1</sup>

with Ross Anderson<sup>2</sup>, Will Ma<sup>4</sup>, Christian Tjandraatmadja<sup>2</sup>, and Juan Pablo Vielma<sup>2,3</sup>

<sup>1</sup>Rice University, Computational and Applied Mathematics
<sup>2</sup>Google Research, Operations Research Team
<sup>3</sup>MIT, Sloan School of Management
<sup>4</sup>Columbia University, Columbia Business School







# Strong mixed-integer programming formulations for trained neural networks

### Joey Huchette<sup>1</sup>

with Ross Anderson<sup>2</sup>, Will Ma<sup>4</sup>, Christian Tjandraatmadja<sup>2</sup>, and Juan Pablo Vielma<sup>2,3</sup> and Yeesian Ng<sup>5</sup> and Ondrej Sykora<sup>2</sup>

<sup>1</sup>Rice University, Computational and Applied Mathematics
<sup>2</sup>Google Research, Operations Research Team
<sup>3</sup>MIT, Sloan School of Management
<sup>4</sup>Columbia University, Columbia Business School
<sup>5</sup>MIT, Operations Research Center







# Strong mixed-integer programming formulations for trained neural networks

### Joey Huchette<sup>1</sup>

with Ross Anderson<sup>2</sup>, Will Ma<sup>4</sup>, Christian Tjandraatmadja<sup>2</sup>, and Juan Pablo Vielma<sup>2,3</sup> and Yeesian Ng<sup>5</sup> and Ondrej Sykora<sup>2</sup> and Craig Boutilier<sup>6</sup>, Martin Mladenov<sup>6</sup>, and Moonkyung Ryu<sup>6</sup>

<sup>1</sup>Rice University, Computational and Applied Mathematics
<sup>2</sup>Google Research, Operations Research Team
<sup>3</sup>MIT, Sloan School of Management
<sup>4</sup>Columbia University, Columbia Business School
<sup>5</sup>MIT, Operations Research Center
<sup>6</sup>Google Research, MudCats



# **The Problem**

### Optimize



End goal is *optimization*: Make the best possible decision  $x^* \in \Omega$ 



### Predict,

Supervised learning: Learn a function using historical input/output data



End goal is *generalization*: Given "reasonable" unseen point  $(x^*, y^*)$ , want  $y^* \approx NN(x^*)$ 



### Optimize



End goal is *optimization*: Make the best possible decision  $x^* \in \Omega$ 



Predict, then Optimize





# Application: Verifying robustness

Adversarial Examples: find small change in the input to change the output



 $\begin{array}{ll} \max & \mathbb{P}[x \text{ is a gibbon}] - \mathbb{P}[x \text{ is a panda}] \\ \text{s.t.} & ||x - (\text{reference panda})||_{\infty} \leq \epsilon \end{array}$ 

#### Predict: Classify input image

Optimize: Prove (or disprove) robustness of model to small perturbations

Very popular research topic in ML community right now

Note that "proofs" of robustness (not just feasible solutions) are very useful!



### Application: Deep reinforcement learning

**Predict:** future costs for action *a* in state *x*:

$$Q(x,a) \approx \mathrm{NN}(x,a)$$

**Optimize:** to pick the lowest cost action:

$$\min_{a} c(x,a) + \alpha \cdot \text{NN}(x,a)$$

In a combinatorial or continuous action space, optimization can be hard!



Teach a cheetah to run (without falling over)



# Application: Designing DNA for protein binding

**Predict:** probability/strength of DNA sequence binding to a given protein

**Optimize:** find the best/many binding sequences (potentially with side constraints)

A typical architecture for predicting protein binding, e.g. (Alipanahi et al 2015, Zeng et al 2016)







## Application: Bayesian optimization

- Goal: Minimize black-box function
  - "Neural architecture search"
  - Drug design
- Sequential experimentation, each sample is costly







Predict: Amount new point improves model

Optimize: Choose best point to sample

- Pros of Gaussian process:
  - Bayesian calculations are straightforward
- (Potential) pros of neural network:
  - Incremental training; handles discontinuities
     + high dims ("model-free"), input constraints
  - Uncertainty models: dropout, ensembling, Thompson sampling, etc.





# The Solution

### How to "Optimize" in "Predict, then Optimize"





# A (not comprehensive) Literature Sampling

- Gradient descent on continuous problems (e.g. deep dream, style transfer, adversarial examples)
- Cross Entropy Method (importance sampling based guesses, brain robotics)
- Input Convex Neural Networks + LP (Amos 2016)
- The big-M MIP (too many authors to list all)
  - Tightening big-Ms: (Tjeng 2017, Fischetti 2018)
- CP/SAT/SMT (often from the "verifiability" community)
  - Reluplex (Katz 2017), Planet (Elhers 2017)
- Lagrangian Relaxation (Dvijotham 2018)
- Abstract Domains/Zonotopes (Singh 2019)
- Nice recent survey (Bunel 2017)



- 1. MIP formulations for popular neural network building blocks Stronger than existing approaches (theoretically and empirically)
  - $\circ$  Most common case: Affine function on box domain  $\rightarrow$  ReLU nonlinearity
  - Other nonlinearities: max-pooling, reduce max, clipped ReLU, ...
  - Other domains: Product of simplices (*one-hot encodings*)
  - Structure: Tight big-M formulation + efficiently separable cuts
- 2. General machinery

Recipes for strong MIP formulations for the maximum of d affine functions

- One for ideal formulations, another for hereditarily sharp ones
- Efficiently separable via subgradient-method
- A series of simplifications under common setting
- 3. tf.opt: Software with TensorFlow-like modeling, multiple backends
- 4. Computational results on verification and reinforcement learning problems



# MIP Formulations for Neural Networks

## MIP formulations in one slide

- A MIP formulation for some set  $S \subseteq \mathbb{R}^n$  is:
  - A polyhedra  $Q \subseteq R^{n+r}$ , where
  - $Proj_x(\{ (x,z) \in Q \mid z \in Z^r \}) = S$
- What makes a MIP formulation good?
  - Size: r is small, Q is "simple"
  - Sharp: Proj<sub>x</sub>(Q) = Conv(S)
  - Hereditarily sharp: Sharp after any fixings of binary variables z
  - Ideal (perfect):  $ext(Q) \subseteq R^n \times Z^r$
- Ideal  $\Rightarrow$  sharp, so...



*Ideal formulations = Best possible = Our goal* 



### Neural networks = Piecewise linear functions

- Standard ReLU-based network:  $x^r = NN(x^0)$ , where  $x_i^j = \max\left\{0, w^{j,i} \cdot x_i^{j-1} + b^{j,i}\right\} \quad \forall j = 1, \dots, r, i = 1, \dots, m_j$
- Big-M formulation for single ReLU neuron

$$egin{cases} L \leq x \leq U \ y = \max\{0, w \cdot x + b\} \ y \geq 0 \ y \geq w \cdot x + b \ y \leq M^+ z \end{cases}$$

$$y \le w \cdot x + b - M^{-}(1 - z)$$
$$(x, y, z) \in [L, U] \times \mathbb{R} \times \{0, 1\}$$



• How strong is it?

is

## MIP formulation strength

• How strong is it? *Not very!* 

Can be arbitrarily bad, even in fixed input dimension



big-M formulation



• How to close the gap?



tightest possible formulation

### An ideal formulation for ReLU neurons

**Theorem** (Anderson, H., Tjandraatmadja, Vielma 2018)

An ideal formulation for  $\{(x, y = \max\{0, w \cdot x + b\}) : L \le x \le U\}$  is 

$$y \ge w \cdot x + b \tag{1a}$$

$$y \le \sum_{i \in I} w_i \left( x_i - L_i (1 - z) \right) + \left( b + \sum_{i \notin I} w_i U_i \right) z \quad \forall I \subseteq [n]$$
(1b)  
$$y_i(z) \in [L, U] \times \mathbb{R}_{\ge 0} \times \{0, 1\}.$$
(1c)

 $(x, y, z) \in [L, U] \times \mathbb{R}_{>0} \times \{0, 1\}.$ 

- Each inequality in (1b) is facet-defining (under very mild conditions).
- Moreover, we can identify the most violated constraint in (1b) in O(n) time.

Big-M formulation = (1a), (1c), and two constraints from (1b)

**Idea:** Start with big-M formulation, use cut callbacks to separate (1b) as-needed



$$\left\{ (x,y) \in D \times \mathbb{R} : y = \max_{k=1}^d w^k \cdot x + b^k \right\}$$

**Step 1:** Write down ideal "multiple choice" formulation (i.e. the "Balas" formulation):

$$egin{aligned} &(x,y) = \sum_{k=1}^d ( ilde{x}^k, w^k \cdot ilde{x}^k + b^k z_k) \ &w^k \cdot ilde{x}^k + b^k z_k \geq w^\ell \cdot ilde{x}^k + b^\ell z_k & orall k, \ell \in [d]: \ k 
eq \ell \ & ilde{x}^k \in z_k \cdot D & orall k \in [d] \ &z \in \Delta^d. \end{aligned}$$



$$\left\{ (x,y) \in D \times \mathbb{R} : y = \max_{k=1}^d w^k \cdot x + b^k \right\}$$

Step 2: Re-write constraints in "set" form:

$$egin{aligned} & (x,y) = \sum_{k=1}^d ( ilde{x}^k, w^k \cdot ilde{x}^k + b^k z_k) \ & ilde{x}^k \in z_k \cdot D_{|k} & orall k \in [d] \ & z \in \Delta^d, \end{aligned}$$

where

$$D_{|k} \equiv \left\{ x \in D : k \in \arg \max_{\ell=1}^{d} w^{\ell} \cdot x + b^{\ell} \right\}$$
$$= \left\{ x \in D : w^{k} \cdot x + b^{k} \ge w^{\ell} \cdot x + b^{\ell} \quad \forall \ell \neq k \right\}$$



$$\left\{ (x,y) \in D \times \mathbb{R} : y = \max_{k=1}^d w^k \cdot x + b^k \right\}$$

Step 3: Rewrite all logic as bounds on output y (a *primal* characterization):

$$\begin{split} y &\leq \overline{g}(x,z) \equiv \max_{\tilde{x}^1,\dots,\tilde{x}^d} \left\{ \sum_{k=1}^d w^k \cdot \tilde{x}^k + b^k z_k : \begin{array}{c} x = \sum_k \tilde{x}^k \\ \tilde{x}^k \in z_k \cdot D_{|k} \quad \forall k \end{array} \right\} \\ y &\geq \underline{g}(x,z) \equiv \min_{\tilde{x}^1,\dots,\tilde{x}^d} \left\{ \sum_{k=1}^d w^k \cdot \tilde{x}^k + b^k z_k : \begin{array}{c} x = \sum_k \tilde{x}^k \\ \tilde{x}^k \in z_k \cdot D_{|k} \quad \forall k \end{array} \right\} \\ x,z) \in D \times \Delta^d, \end{split}$$

where

$$D_{|k} \equiv \left\{ x \in D : k \in \arg \max_{\ell=1}^{d} w^{\ell} \cdot x + b^{\ell} \right\}$$
$$= \left\{ x \in D : w^{k} \cdot x + b^{k} \ge w^{\ell} \cdot x + b^{\ell} \quad \forall \ell \neq k \right\}$$



$$\left\{ (x,y) \in D \times \mathbb{R} : y = \max_{k=1}^d w^k \cdot x + b^k \right\}$$

Step 4: Apply Lagrangian relaxation to aggregation constraints (a dual characterization)

$$y \leq \overline{\alpha} \cdot x + \sum_{k=1}^{d} \left( \max_{x^k \in D_{|k}} \{ (w^k - \overline{\alpha}) \cdot x^k \} + b^k \right) z_k \quad \forall \overline{\alpha} \in \mathbb{R}^{\eta}$$
$$y \geq \underline{\alpha} \cdot x + \sum_{k=1}^{d} \left( \min_{x^k \in D_{|k}} \{ (w^k - \underline{\alpha}) \cdot x^k \} + b^k \right) z_k \quad \forall \underline{\alpha} \in \mathbb{R}^{\eta}$$
$$(x, z) \in D \times \Delta^d.$$

End of analysis: in general,  $D_{|k|}$  is complicated.

Can separate over via subgradient method, or...



$$\left\{ (x,y) \in D \times \mathbb{R} : y = \max_{k=1}^d w^k \cdot x + b^k \right\}$$

**Step 1:** Write down ideal "multiple choice" formulation (i.e. the "Balas" formulation):

$$egin{aligned} &(x,y) = \sum_{k=1}^d ( ilde{x}^k, w^k \cdot ilde{x}^k + b^k z_k) \ &w^k \cdot ilde{x}^k + b^k z_k \geq w^\ell \cdot ilde{x}^k + b^\ell z_k & orall k, \ell \in [d]: \ k 
eq \ell \ & ilde{x}^k \in z_k \cdot D & orall k \in [d] \ &z \in \Delta^d. \end{aligned}$$



$$\left\{ (x,y) \in D \times \mathbb{R} : y = \max_{k=1}^d w^k \cdot x + b^k \right\}$$

Step 2: Re-write constraints in "set" form:

$$egin{aligned} & (x,y) = \sum_{k=1}^d ( ilde{x}^k, w^k \cdot ilde{x}^k + b^k z_k) \ & ilde{x}^k \in z_k \cdot D_{|k} & orall k \in [d] \ & z \in \Delta^d, \end{aligned}$$

where

$$D_{|k} \equiv \left\{ x \in D : k \in \arg \max_{\ell=1}^{d} w^{\ell} \cdot x + b^{\ell} \right\}$$
$$= \left\{ x \in D : w^{k} \cdot x + b^{k} \ge w^{\ell} \cdot x + b^{\ell} \quad \forall \ell \neq k \right\}$$



$$\left\{ (x,y) \in D \times \mathbb{R} : y = \max_{k=1}^d w^k \cdot x + b^k \right\}$$

**Step 3:** Relax domain constraint:

$$egin{aligned} &(x,y) = \sum_{k=1}^d ( ilde{x}^k, w^k \cdot ilde{x}^k + b^k z_k) \ & ilde{x}^k \in z_k \cdot D & orall k \in [d] \ &z \in \Delta^d. \end{aligned}$$

Still a valid formulation.



$$\left\{ (x,y) \in D \times \mathbb{R} : y = \max_{k=1}^d w^k \cdot x + b^k \right\}$$

Step 4: Rewrite all logic as bounds on output y (a *primal* characterization):

$$\begin{split} y &\leq \overline{h}(x,z) \equiv \max_{\tilde{x}^1,\dots,\tilde{x}^d} \left\{ \sum_{k=1}^d w^k \cdot \tilde{x}^k + b^k z_k : \begin{array}{c} x = \sum_k \tilde{x}^k \\ \tilde{x}^k \in z_k \cdot D \quad \forall k \end{array} \right\} \\ y &\geq \underline{h}(x,z) \equiv \min_{\tilde{x}^1,\dots,\tilde{x}^d} \left\{ \sum_{k=1}^d w^k \cdot \tilde{x}^k + b^k z_k : \begin{array}{c} x = \sum_k \tilde{x}^k \\ \tilde{x}^k \in z_k \cdot D \quad \forall k \end{array} \right\} \\ x,z) \in D \times \Delta^d. \end{split}$$



$$\left\{ (x,y) \in D \times \mathbb{R} : y = \max_{k=1}^d w^k \cdot x + b^k \right\}$$

**Step 5:** Replace lower bounds for a *hereditarily sharp* formulation:

$$\begin{split} y &\leq \overline{h}(x,z) \equiv \max_{\tilde{x}^1, \dots, \tilde{x}^d} \left\{ \sum_{k=1}^d w^k \cdot \tilde{x}^k + b^k z_k : \begin{array}{l} x = \sum_k \tilde{x}^k \\ \tilde{x}^k \in z_k \cdot D \quad \forall k \end{array} \right\} \\ y &\geq w^k \cdot x + b^k \quad \forall k \in [d] \\ (x,z) \in D \times \Delta^d. \end{split}$$



$$\left\{ (x,y) \in D \times \mathbb{R} : y = \max_{k=1}^d w^k \cdot x + b^k \right\}$$

**Step 6:** Apply Lagrangian relaxation to aggregation constraints:

$$\begin{split} y &\leq \overline{\alpha} \cdot x + \sum_{k=1}^{d} \left( \max_{x^k \in D} \{ (w^k - \overline{\alpha}) \cdot x^k \} + b^k \right) z_k \qquad \quad \forall \overline{\alpha} \in \mathbb{R}^{\eta} \\ y &\geq w^k \cdot x + b^k \quad \forall k \in [d] \\ (x, z) \in D \times \Delta^d. \end{split}$$



$$\left\{ (x,y) \in D \times \mathbb{R} : y = \max_{k=1}^d w^k \cdot x + b^k \right\}$$

**Step 7:** Analyze further:

- **Proposition** If d=2, then the hereditarily sharp formulation is ideal.
- **Proposition** If the domain is a product of simplices, separation over hereditarily sharp formulation reduces to a transportation problem.
- **Proposition** If d=2 and the domain is a product of simplices, the transportation problem has a closed form solution and efficient separation.



### An ideal formulation for ReLU neurons

**Theorem** (Anderson, H., Tjandraatmadja, Vielma 2018)

An ideal formulation for  $\{(x, y = \max\{0, w \cdot x + b\}) : L \le x \le U\}$  is 

$$y \ge w \cdot x + b \tag{1a}$$

$$y \le \sum_{i \in I} w_i \left( x_i - L_i (1 - z) \right) + \left( b + \sum_{i \notin I} w_i U_i \right) z \quad \forall I \subseteq [n]$$
(1b)  
$$y_i(z) \in [L, U] \times \mathbb{R}_{\ge 0} \times \{0, 1\}.$$
(1c)

 $(x, y, z) \in [L, U] \times \mathbb{R}_{>0} \times \{0, 1\}.$ 

- Each inequality in (1b) is facet-defining (under very mild conditions).
- Moreover, we can identify the most violated constraint in (1b) in O(n) time.

Big-M formulation = (1a), (1c), and two constraints from (1b)

**Idea:** Start with big-M formulation, use cut callbacks to separate (1b) as-needed



# Hereditarily sharp formulation for max-of-d affine functions

- Max-of-d affine functions ≡
  - Max pooling (small d)
  - Reduce max (large d)

12	20	30	0			
8	12	2	0	$2 \times 2$ Max-Pool	20	30
34	70	37	4		112	37
112	100	25	12			

Proposition (Anderson, H., Ma, Tjandraatmadja, Vielma 2019)

• A hereditarily sharp MIP formulation for  $\left\{ \left(x, y = \max\left\{w^k \cdot x + b^k\right\}_{k=1}^d\right) : L \le x \le U \right\}$  is:

$$\begin{split} y &\leq \sum_{i=1}^{\eta} \left( w_i^{I(i)} x_i + \sum_{k=1}^d \max\{(w_i^k - w_i^{I(i)}) L_i, (w_i^k - w^{I(i)}) U_i\} z_k \right) + \sum_{k=1}^d b^k z_k \quad \forall I : [\eta] \to [d] \\ y &\geq w^k \cdot x + b^k \quad \forall k \in [d] \\ x, z) \in D \times \Delta^d \\ z \in \{0, 1\}^d. \end{split}$$

• The most violated inequality can be identified in time linear in O(dn) time.



# **Computational Results**

### Network 1: Small network with standard training



### Network 2: Small network with L1 regularization (Xiao et al. 2019)



# Q-learning: Are optimal actions tractable?

MIP Action Selection Performance h256-128





Ongoing joint work with Craig Boutilier, Martin Mladenov, and Moonkyung Ryu (Google Research)

### Q-learning: Are optimal actions better?

#### Cheetah Reward by Policy Iteration



- Gradient Descent h32-16 - MIP h32-16

**Policy Iterations** 

Ongoing joint work with Craig Boutilier, Martin Mladenov, and Moonkyung Ryu (Google Research)



### Q-learning: Are optimal actions better?

#### Cheetah Reward by Policy Iteration



**Policy Iterations** 

Ongoing joint work with Craig Boutilier, Martin Mladenov, and Moonkyung Ryu (Google Research)



### Conclusion

- Strong MIP formulations for optimizing over trained neural networks
- Applications abound: verification, drug design, reinforcement learning, etc.
- Framework of independent interest:

Recipes for strong formulations for the max of d affine functions

- Questions going forward:
  - Separation-based algorithm and implementation
  - How to train the network for optimization in a *principled* way
  - How best to formulate entire network (not just individual neurons)

