# Active Learning of Deep Neural Networks via Gradient-Free Cutting Planes

Erica Zhang[*], Fangzhao Zhang[*], Mert Pilanci

Stanford | ENGINEERING

June 16, 2025

# Introduction

# Active Learning of Deep Networks

**Active Learning (AL)** is a data acquisition paradigm for supervised learning. Active learning aims at training a model $\mathcal{X} \to \mathcal{Y} : f(x; \theta) = y$, characterized by its parameters $\theta$, with the *most informative* data $\subseteq \mathcal{X} \times \mathcal{Y}$.

▶ As deep networks scale, training efficiency increasingly hinges on selecting data that contributes the most to learning, rather than relying on random samples.

  • A key class of active learning algorithms selects the next point $(x, y)_{k+1}$ to *maximize model uncertainty reduction*, leveraging the model trained on the current labeled set $\{(x_i, y_i)\}_{i=1}^{k}$.

▶ However, the *nonconvexity* of deep neural networks poses challenges for active learning.

  •
    **Challenge 1:** Existing deep AL methods lack convergence guarantees and many rely heavily on heuristics.

**Cutting-plane methods** are first introduced by Gomory, 1958 [1] for LPs and now heavily used in commercial MILP solvers.
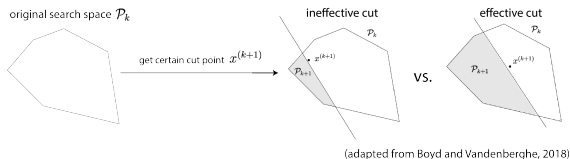


Figure 2: The Cutting-plane method

▶ Consider any minimization problem with objective $f(\theta)$ whose solution set $\Theta$ is a convex set.

▶ Assuming existence of an oracle such that given any input $\theta_0$, we receive a cutting plane that cuts between the current input $\theta_0$ and the desired solution set $\Theta$.

▶ Such cut is given by subgradient of $f$ at query point $\theta_0$.

▶ **The linear cutting-plane AL.** Louche & Ralaivola, 2015 [2] proposes the first use of cutting-plane method in the context of active learning.

- **Challenge 2:** Restricted to AL for shallow (linear) models.
- Allows only linear cuts and binary classifications.
- *Our proposed Cutting-Plane Active Learning (CPAL) method addresses these challenges while preserving the theoretical convergence guarantees inherent to cutting-plane methods.*

---

**Algorithm 3** Generic (Linear) Cutting-Plane AL

1: $\mathcal{T}^0 \leftarrow \mathcal{B}_2$
2: $t \leftarrow 0$
3: **repeat**
4:     $\theta_c^t \leftarrow \text{center}(\mathcal{T}^t)$
5:     $x_{n_t}, y_{n_t} \leftarrow \text{QUERY}(\mathcal{T}^t, \mathcal{D})$
6:     **if** $y_{n_t}\langle \theta_c^t, x_{n_t}\rangle < 0$ **then**
7:        $\mathcal{T}^{t+1} \leftarrow \mathcal{T}^t \cap \{z : y_{n_t}\langle z, x_{n_t}\rangle \geq 0\}$
8:        $t \leftarrow t + 1$
9:     **end if**
10: **until** $\mathcal{T}^t$ is small enough
11: **return** $\theta_c^t$

1: **function** QUERY($\mathcal{T}, \mathcal{D}$)
2:     Sample $M$ points $s_1, \ldots, s_M$ from $\mathcal{T}$
3:     $g \leftarrow \frac{1}{M}\sum_{k=1}^{M} s_k$
4:     $x \leftarrow \arg\min_{x_i \in \mathcal{D}}\langle g, x_i\rangle$
5:     $y \leftarrow$ get label from an expert
6:     **return** $x, y$
7: **end function**

# The Cutting-Plane Active Learning (CPAL) Method

- **The Cutting-Plane Active Learning (CPAL) method** addresses the **two** aforementioned key challenges:
  - **Theorem 4.2 (linear to non-linear):** Reformulates deep ReLU network training as a linear program via activation patterns.
    - Extends to *multi-class classification* (via one-vs-rest decomposition) and *regression* (by adapting constraints to bound prediction error).
    - **Contribution 1:** Generalizes beyond the linearity of traditional cutting-plane methods to nonlinear models.
  - **Theorem 6.3 (convergence guarantees):** Establishes volumetric convergence guarantees in parameter space.
    - Preserves the classic $(1 - 1/e)$ volumetric shrinkage guarantee of cutting-plane methods with cuts through center of gravity [3].
    - Two-layer NN prediction functions converge in norm to the optimal decision boundary under CPAL (Corollary 6.4).
    - **Contribution 2:** Presents the first deep AL method with convergence guarantees.

# The CPAL Algorithm

**Algorithm 4** Cutting-plane AL for Binary Classification with Limited Queries

1: $\mathcal{T}^0 \leftarrow \mathcal{B}_2$
2: $t \leftarrow 0$
3: $\mathcal{D}_{\mathrm{AL}} \leftarrow \mathbf{0}$
4: **repeat**
5:      $\theta_c^t \leftarrow \mathrm{center}(\mathcal{T}^t)$
6:      **for** $s$ in $\{1, -1\}$ **do**
7:          $(x_{n_t}, y_{n_t}) \leftarrow \mathrm{QUERY}(\mathcal{T}^t, \mathcal{D} \setminus \mathcal{D}_{\mathrm{AL}}, s)$
8:          **if** $y_{n_t} \cdot f^{\mathrm{two\text{-}layer}}(x_{n_t}; \theta_c^t) < 0$ **then**
9:              $\mathcal{D}_{\mathrm{AL}} \leftarrow \mathrm{ADD}(\mathcal{D}_{\mathrm{AL}}, (x_{n_t}, y_{n_t}))$
10:            $\mathcal{T}^{t+1} \leftarrow \mathcal{T}^t \cap \{\theta : y_{n_t} \cdot f^{\mathrm{two\text{-}layer}}(x_{n_t}; \theta) \geq 0, \mathcal{C}(\{n_t\}), \mathcal{C}'(\{n_t\})\}$
11:            $t \leftarrow t + 1$
12:          **end if**
13:      **end for**
14: **until** $|\mathcal{D}_{\mathrm{AL}}| \geq n_{\mathrm{budget}}$
15: **return** $\theta_c^t$

1: **function** QUERY($\theta, s$)
2:      $(x, y) \leftarrow \arg\min_{(x_i, y_i) \in \mathcal{D}_{\mathrm{QS}}} s f^{\mathrm{two\text{-}layer}}(x_{n_t}; \theta)$
3:      **return** $(x, y)$
4: **end function**

# Theorem 4.2 (Linear to Non-Linear)

Consider a ReLU network with $n$ hidden layers for binary classification:

$$\text{find} \qquad W_1, W_2, \ldots, W_{n+1}$$
$$\text{s.t.} \qquad y \odot (\cdots(((XW_1)_+ W_2)_+ \cdots)_+ W_n)_+ W_{n+1} \geq 1. \tag{1}$$

Using the notations established in the paper, then when $m_i \geq \prod_{i=1}^n 2P_i$ for each $i \in [n]$, Problem (1) is equivalent to:

$$\text{find} \qquad u_{j_n j_{n-1} \ldots j_1}^{c_n c_{n-1} \ldots c_1}$$

$$\text{s.t.} \qquad y \odot \sum_{j_n=1}^{P_n} D_{j_n}^{(n)} \left( \mathcal{T}_1^{(n-1)}(D^{(n-1)}) - \mathcal{T}_2^{(n-1)}(D^{(n-1)}) \right) \geq 1$$

$$(2D_{j_i}^{(i)} - I) \mathcal{T}_{c_{n-1} \ldots c_{i-1}}^{(n-1) \ldots (i-1)}(D^{(i-1)}) \geq 0, \quad 2 \leq i \leq n \tag{2}$$

$$(2D_{j_1}^{(1)} - I) X u_{j_n \ldots j_1}^{c_n \ldots c_1} \geq 0.$$

where $c_i \in \{1, 2\}$, and

$$\mathcal{T}_{c_{n-1} \ldots c_i}^{(n-1) \ldots (i)}(D^{(i)}) = \sum_{j_i=1}^{P_i} D_{j_i}^{(i)} \left( \mathcal{T}_{c_{n-1} \ldots c_i 1}^{(n-1) \ldots (i)(i-1)}(D^{(i-1)}) \right.$$

$$\left. - \mathcal{T}_{c_{n-1} \ldots c_i 2}^{(n-1) \ldots (i)(i-1)}(D^{(i-1)}) \right), \quad \forall i \leq n-1$$

$$\mathcal{T}_{c_{n-1} \ldots c_1}^{(n-1) \ldots (1)}(D^{(1)}) = \sum_{j_1=1}^{P_1} D_{j_1}^{(1)} X \left( u_{j_n \ldots j_1}^{1 c_{n-1} \ldots c_1} - u_{j_n \ldots j_1}^{2 c_{n-1} \ldots c_1} \right).$$

**Theorem 6.3 (Convergence via Center of Gravity)**

Let $\mathcal{T} \subseteq \mathbb{R}^d$ be a convex body and let $\theta_G$ denote its center of gravity. The polyhedron cut given in Algorithm 4 (CPAL), i.e.,

$$\mathcal{T} \cap \left\{ \theta : y_n \cdot f^{\text{two-layer}}(x_n; \theta) \geq 0, \ \mathcal{C}(\{n\}), \ \mathcal{C}'(\{n\}) \right\},$$

where the pair $(x_n, y_n)$ is the data point returned by the cutting-plane oracle after receiving the query point $\theta_G$, partitions $\mathcal{T}$ into two subsets:
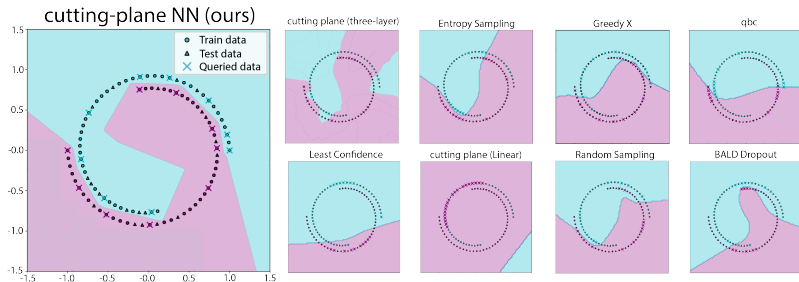
$$\mathcal{T}_1 := \left\{ \theta \in \mathcal{T} : y_n \cdot f^{\text{two-layer}}(x_n; \theta) \geq 0, \ \mathcal{C}(\{n\}), \ \mathcal{C}'(\{n\}) \right\},$$

$$\mathcal{T}_2 := \left\{ \theta \in \mathcal{T} : y_n \cdot f^{\text{two-layer}}(x_n; \theta) < 0 \ \vee \ \neg\mathcal{C}(\{n\}) \ \vee \ \neg\mathcal{C}'(\{n\}) \right\},$$

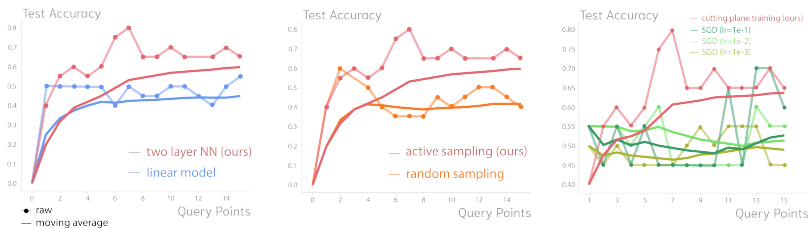where $\neg$ denotes the complement. Then $\mathcal{T}_1$ satisfies:

$$\mathbf{vol}(\mathcal{T}_1) < \left( 1 - \frac{1}{e} \right) \cdot \mathbf{vol}(\mathcal{T}).$$

# Experiments

Figure 3: Decision boundaries for binary classification on the spiral dataset for the cutting-plane AL method using a two-layer ReLU neural network, alongside various deep AL baselines. For compactness, we also include the decision boundaries for the cutting-plane AL method with a three-layer ReLU network in the collage to demonstrate its feasibility. For fairness of comparison, we use the same two-layer ReLU network structure and embedding size of 623 for all methods. We enforce the same hyperparameters for all deep AL baselines and select the best performing number of training epochs at 2000 and a learning rate at 0.001 to ensure optimal performance.

Figure 4: Sentiment analysis on IMDB movie review dataset with two-layer ReLU model. We take Phi-2 embedding as our training features and compare with various baselines. The result shows that the introduction of non-linearity improves upon linear model performance, our active sampling scheme effectively identifies valuable training points compared to random sampling, and our cutting-plane training scheme is more effective than SGD in this setting. Linear and our reframed two-layer models are initialized to predict zero while two-layer NN trained with SGD has random weight initialization, thus starting from non-zero prediction.

# Summary

- We generalize classic cutting-plane methods from linear models to nonlinear deep ReLU networks via activation pattern enumeration.

- We propose the first deep active learning algorithm with formal convergence guarantees.

- We establish the viability of gradient-free training through cutting-plane techniques.

- Our cutting-plane active learning method (CPAL) outperforms popular deep AL baselines on both synthetic and real-world datasets.

- While GPU-based CPAL currently faces scalability limitations, it lays a theoretical foundation for deep AL and offers a promising direction as LP solvers and activation pattern sampling continue to improve.

# Thank you!

# Bibliography

[1] R. E. Gomory, "An algorithm for integer solutions to linear programs,", 1958. [Online]. Available: https://api.semanticscholar.org/CorpusID:116324171.

[2] L. Ralaivola and U. Louche, *From cutting planes algorithms to compression schemes and active learning*, 2015. arXiv: 1508.02986 [cs.LG]. [Online]. Available: https://arxiv.org/abs/1508.02986.

[3] B. Grünbaum, "Partitions of mass-distributions and of convex bodies by hyperplanes," *Journal of Mathematical Analysis and Applications*, vol. 10, no. 4, pp. 1257–1261, 1960, Received January 22, 1960. This research was supported by the United States Air Force through the Air Force Office of Scientific Research of the Air Research and Development Command, under contract No. AF49(638)-253.