CROSS-ENTROPY LOSS OF APPROXIMATED DEEP NEURAL NETWORKS



JIŘÍ ŠÍMA & PETRA VIDNEROVÁ



Efficient Processsing of Deep Neural Networks

(Sze, Chen, Yang, Emer, Morgan & Claypool Publishers, 2020)

The energy efficiency of DNNs on low-power, battery-operated embedded hardware (e.g., cellphones, smartwatches, augmented reality glasses) is a critical challenge.

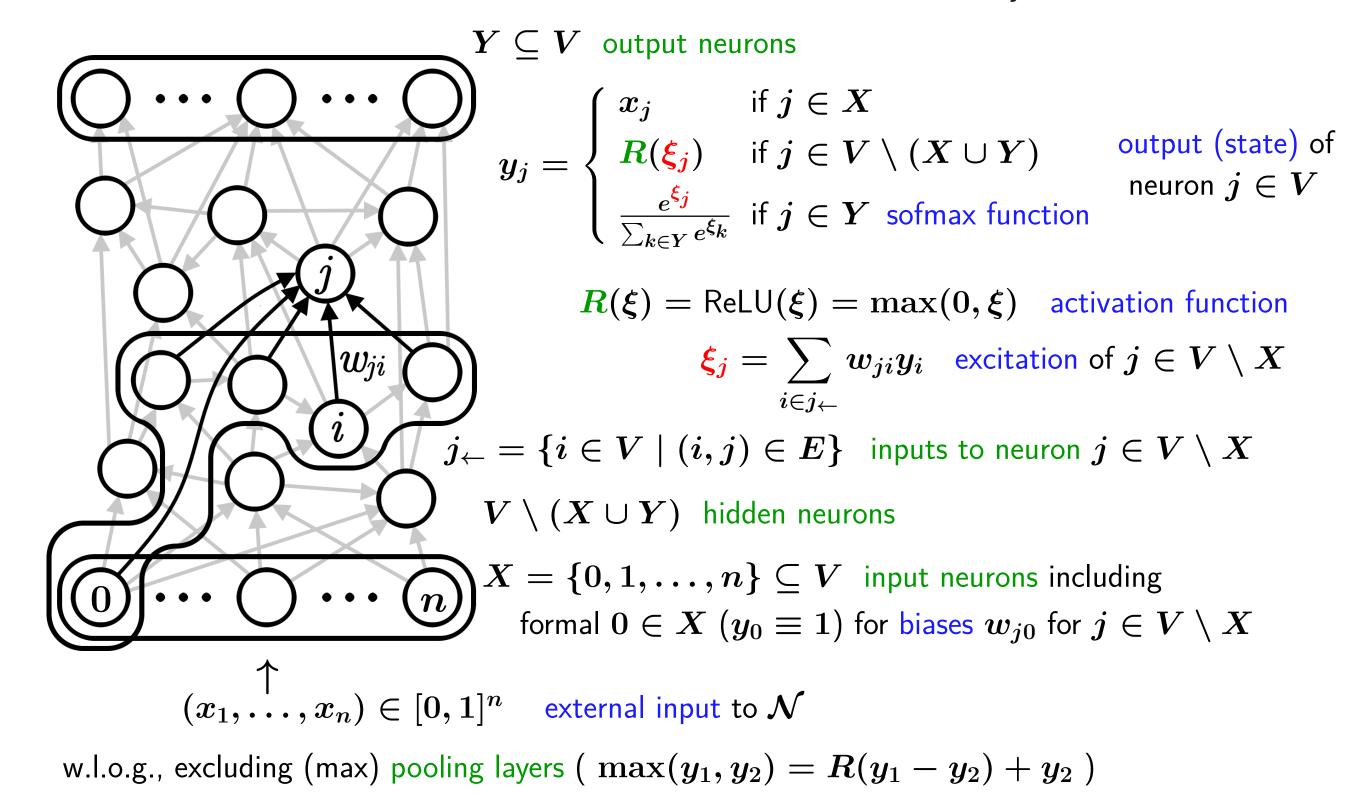
---- reducing the energy cost of DNNs:

- 1. Specialized Hardware Accelerators for DNN inference (on GPUs, FPGAs, in-memory, etc.)
- 2. Approximate Computing in error-tolerant applications (e.g. image classification) save large amounts of energy with minimal accuracy loss by reducing:
- Model Size: pruning, compression, weight sharing, approximate multipliers
- Arithmetic Precision: fixed-point operations, reduced weight bit-width, nonuniform quantization

The aim of this study: Estimate the maximum cross-entropy loss of approximated classification DNNs.

Formal Model of DNN

The architecture of a DNN \mathcal{N} is a connected directed acyclic graph (V, E) composed of neurons, where edges $(i,j) \in E \subset V \times V$ are labeled with real-valued weights w_{ii} :



1. Tool: Shortcut Weights

The excitation ξ_i of any neuron $j \in V \setminus X$ is a continuous piecewise linear function of the external input

$$\rightarrow \quad \xi_j = \sum_{i \in X} W_{ji} y_i \quad \text{for } (y_1, \dots, y_n) \in \Xi$$

within a subset $\Xi \subset [0,1]^n$ of the input space, where W_{ii} are referred to as the shortcut weights (bias) from input neurons $i \in X$ to neuron $j \in V \setminus X$.

For an input $(x_1, \ldots, x_n) \in [0, 1]^n$, let $\Xi_S \subset [0, 1]^n$ be its neighborhood within which ξ_i are linear for all $j \in V \setminus X$ under fixed shortcut weights, where

$$S = S(x_1, \ldots, x_n) = \{j \in V \setminus (X \cup Y) \mid \xi_i < 0\}$$

denotes the set of hidden neurons saturated at zero output, $y_i = R(\xi_i) = 0$.

 $\rightarrow \Xi_S$ is a convex polytope—an intersection of finitely many half-spaces:

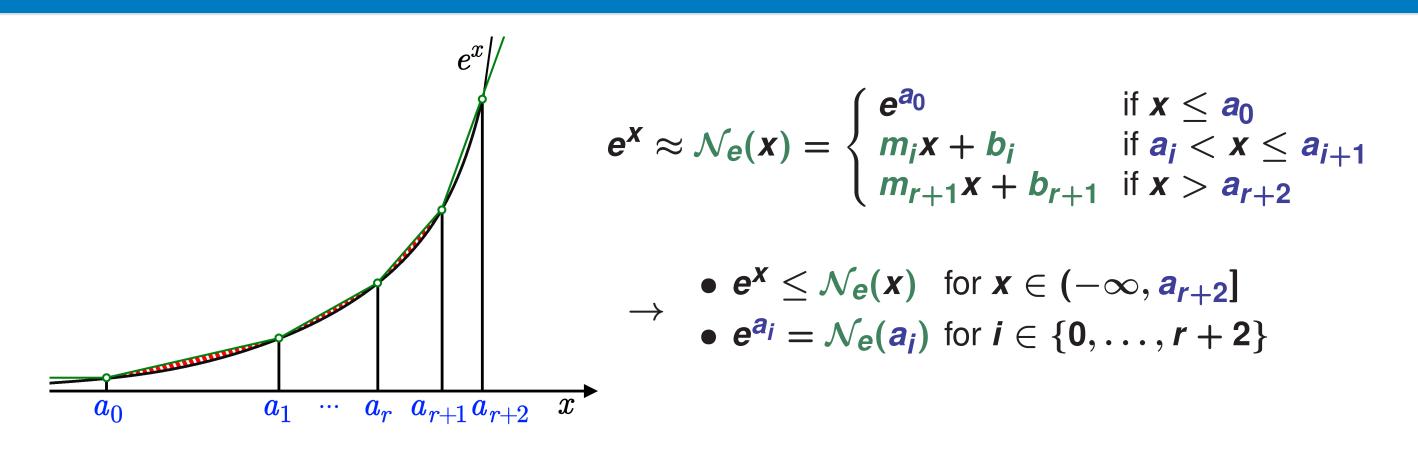
$$\xi_j = \sum_{i \in X} W_{ji} y_i \begin{cases} < 0 \text{ if } j \in S \\ \ge 0 \text{ if } j \notin S \end{cases} \text{ for } j \in V \setminus (X \cup Y)$$

 $0 \le y_i \le 1$ for $i \in X$

Efficient Computation of the shortcut weights via feedforward propagation (skipping S) for $j \in V \setminus X$:

$$W_{ji} = \sum_{k \in j_{\leftarrow} \setminus S} w_{jk} \ W_{ki}$$
 for all $i \in X$ (initially $W_{ki} = \begin{cases} 1 & \text{if } k = i \\ 0 & \text{otherwise} \end{cases}$ for $k, i \in X$)

2. Tool: Continuous Piecewise Linear Interpolation of e^x



Theorem: There are r unique points a_1, \ldots, a_r inside $[a_0, a_{r+1}]$ that minimize

$$\sum_{i=0}^{r} \int_{a_{i}}^{a_{i+1}} \left(m_{i}x + b_{i} - e^{x} \right) dx \qquad \left(\rightarrow e^{a_{i}} = \frac{e^{a_{i+1}} - e^{a_{i-1}}}{a_{i+1} - a_{i-1}} \right).$$

Evaluating the linear interpolation using ReLU networks:

$$\mathcal{N}_{e}(x) = e^{a_0} + \sum_{i=0}^{r+1} R(m_i x + b_i - e^{a_i}) - \sum_{i=0}^{r} R(m_i x + b_i - e^{a_{i+1}})$$

Approximated DNNs

 \mathcal{N} is an approximation of \mathcal{N} , sharing the same input neurons ($\mathbf{X} = \mathbf{X}$) and the same number of output neurons ($|\mathbf{Y}| = |\mathbf{Y}|$) (tilde denotes parameters of \mathcal{N}),

e.g., the weights of \mathcal{N} are rounded to a given number of binary digits in their floating-point representations

Cross-Entropy Loss of Approximated Classification DNNs

The classification error of $\widetilde{\mathcal{N}}$ for an input $(x_1,\ldots,x_n)\in[0,1]^n$, measured by the cross-entropy loss between the softmax categorical probability distributions of $\mathcal N$ and $\mathcal N$:

$$L(x_1, \ldots, x_n) = -\sum_{k \in Y} y_k \ln \widetilde{y_k}$$
 (upper bounds the Kullback-Leibler divergence of \mathcal{N} from $\widetilde{\mathcal{N}}$)

Theorem: It is NP-hard to find the maximum of the cross-entropy loss L over the input space $[0,1]^n$.

Upper-Bounding L Using ReLU Networks $\widetilde{\mathcal{N}_{ck}}$ & \mathcal{N}_c

within the convex polytope Ξ_S around a data point $(x_1, \ldots, x_n) \in T$ of category $c \in Y$ from a test/training set $T \subset [0,1]^n$, where $S = S(x_1, \ldots, x_n)$, restricted to inputs in Ξ_S that are classified by \mathcal{N} into the category \boldsymbol{c} with probability at least \boldsymbol{p} (e.g., $\boldsymbol{p}=0.8$): $\boldsymbol{y_c}>\boldsymbol{p}$

$$\sum_{j \in Y \setminus \{c\}} e^{\xi_j - \xi_c} \stackrel{\xi_j - \xi_c \le 0 \le a_{r+2}}{\le} \sum_{j \in Y \setminus \{c\}} \mathcal{N}_e\left(\xi_j - \xi_c\right) = \mathcal{N}_c(x_1, \dots, x_n) \le \frac{1}{p} - 1 \quad \stackrel{\text{softmax}}{\longrightarrow} y_c \ge p$$

$$L(X_{1},\ldots,X_{n}) = \sum_{k\in Y} y_{k} \ln \frac{1}{\widetilde{y_{k}}} \leq y_{c} \ln \frac{1}{\widetilde{y_{c}}} + (1-y_{c}) \max_{k\in \widetilde{Y}\setminus\{c\}} \ln \frac{1}{\widetilde{y_{k}}}$$

$$\sum_{k\in \widetilde{Y}\setminus\{c\}} \widetilde{\xi_{k}} - \widetilde{\xi_{k}} + y_{c}(\widetilde{\xi_{k}} - \widetilde{\xi_{c}}) \xrightarrow{y_{c} \geq p} \lim_{k \to \infty} \sum_{k\in \widetilde{Y}\setminus\{c\}} \widetilde{\xi_{k}} + R(\widetilde{\xi_{k}} - \widetilde{\xi_{c}}) - pR(\widetilde{\xi_{c}} - \widetilde{\xi_{k}})$$

$$= \max_{k \in \widetilde{Y} \setminus \{c\}} \ln \sum_{j \in \widetilde{Y}} e^{\widetilde{\xi}_{j} - \widetilde{\xi}_{k} + y_{c}\left(\widetilde{\xi}_{k} - \widetilde{\xi}_{c}\right)} \stackrel{y_{c} \geq p}{\leq} \ln \max_{k \in \widetilde{Y} \setminus \{c\}} \sum_{j \in \widetilde{Y}} e^{\widetilde{\xi}_{j} - \widetilde{\xi}_{k} + R\left(\widetilde{\xi}_{k} - \widetilde{\xi}_{c}\right) - pR\left(\widetilde{\xi}_{c} - \widetilde{\xi}_{k}\right)}$$

$$(*) \qquad (*) \qquad (*)$$

$$\overset{(\star)}{\leq} \ln \max_{k \in \widetilde{Y} \setminus \{c\}} \sum_{j \in \widetilde{Y}} \mathcal{N}_{e} \left(\widetilde{\xi_{j}} - \widetilde{\xi_{k}} + R \left(\widetilde{\xi_{k}} - \widetilde{\xi_{c}} \right) - pR \left(\widetilde{\xi_{c}} - \widetilde{\xi_{k}} \right) \right) = \ln \max_{k \in \widetilde{Y} \setminus \{c\}} \widetilde{\mathcal{N}_{ck}}(x_{1}, \dots, x_{n})$$

AppMax Method

Input: DNN \mathcal{N} , its approximation $\overline{\mathcal{N}}$, $(x_1, \ldots, x_n) \in T$ of category $c \in Y$ Output: an upper bound on

$$\max_{(y_1,...,y_n)\in\bigcup_{k\in\widetilde{Y}\setminus\{c\}}\Xi_k^*} L(y_1,\ldots,y_n)$$

Algorithm: For each $k \in Y \setminus \{c\}$ do

- Compose \mathcal{N}^* from \mathcal{N}_{ck} & \mathcal{N}_{c} .
- Determine the saturated neurons $S^* = S^*(x_1, \dots, x_n)$ in \mathcal{N}^* .
- Compute the shortcut weights W_{ii}^* of \mathcal{N}^* for all $j \in V^* \setminus X^*$ and $i \in X^*$.
- Solve the linear program (LP) to find (y_1, \ldots, y_n) that

maximize
$$\widetilde{\mathcal{N}}_{ck}(y_1,\ldots,y_n)$$
 $(\to \sigma_k)$ over the polytope $\Xi_k^*\subseteq \overline{\Xi_{S^*}}$ defined by $\xi_j^*=\sum_{i\in X}W_{ji}^*\,y_i\left\{ egin{array}{l} \le 0 & \text{if } j\in S^* \\ \ge 0 & \text{if } j\notin S^* \end{array} \right.$ for $j\in V^*\setminus (X^*\cup Y^*)$, $\mathcal{N}_c(x_1,\ldots,x_n)\le \frac{1}{p}-1 \qquad \left(\text{where } e^{a_0}\le \frac{1}{p}-1\le e^{a_{r+2}}\right),$ $\widetilde{\xi}_j-\widetilde{\xi}_k+R\left(\widetilde{\xi}_k-\widetilde{\xi}_c\right)-pR\left(\widetilde{\xi}_c-\widetilde{\xi}_k\right)\le a_{r+2} \quad (\star), \qquad (y_1,\ldots,y_n)\in [0,1]^n.$

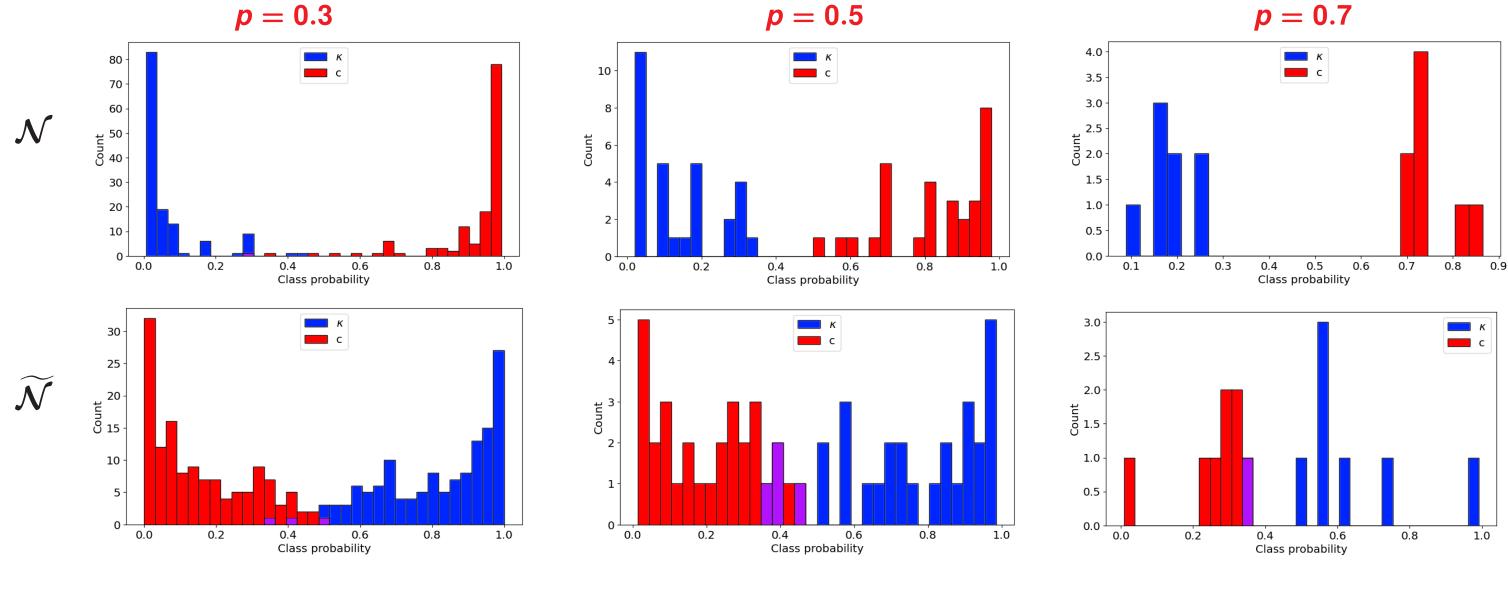
and output $\ln \max_{\pmb{k} \in \widetilde{\pmb{Y}} \setminus \{\pmb{c}\}} \sigma_{\pmb{k}}$

Experiments

- DNN:
 ✓ with 3 fully connected layers (784–64–32–10) and 32-bit weights, trained on the MNIST dataset; approximated as N with weights rounded to 4 bits
- Source Code: https://github.com/PetraVidnerova/ClassificationRoundingErrors
- Test Set: T with 1135 data points of class c = 1, on which \mathcal{N} and \mathcal{N} achieve 100% and 99.91% accuracy, respectively
- Linear Interpolation of e^x : r = 14 optimal points & $a_0 = -5$, $a_{r+1} = 5$, $a_{r+2} = 20$

Number of "misclassified polytopes" around data points in *T*, including so-called misclassified inputs with the maximum estimated upper bound of cross-entropy loss, which are classified correctly as c = 1 by \mathcal{N} with probability at least p, but misclassified as $\kappa \neq c$ by \mathcal{N} :

Histograms of probabilities y_c and y_k over the misclassified inputs:



Future Research Directions

- Broaden AppMax evaluation to other datasets (e.g., CIFAR-100, ImageNet).
- Most suitable error metric (KL divergence, cross-entropy loss, accuracy) with optimal upper bound.
- AppMax for classification DNNs with softmax, via nonlinear Karush-Kuhn-Tucker optimization.
- Approximate global error by estimating the probabilities of convex polytopes from their volumes, measured using the average mean width (evaluated by LP).
- Identify DNN components that can be neglected (e.g., specific weights to be rounded) under explicitly bounded output error.