PAC-Bayesian Learning of Aggregated Binary Activated Neural Networks with Probabilities over Representations

Benjamin Leblanc Louis Fortier-Dubois Gaël Letarte Pascal Germain François Laviolette

June 6th, 2023

June 6th, 2023

• Neural networks: why so difficult to analyse?







→

- Neural networks: why so difficult to analyse?
- Simplification as a solution







< ∃ ▶

- Neural networks: why so difficult to analyse?
- Simplification as a solution
- Probability distribution over the weights of the model







2/23

- Neural networks: why so difficult to analyse?
- Simplification as a solution
- Probability distribution over the weights of the model
- Based on the work of Letarte et al.[1]

[1] Gaël Letarte, Pascal Germain, Benjamin Guedj, and François Laviolette. Dichotomize and generalize: PAC-Bayesian binary activated deep neural networks. NeurIPS, 2019.







• PAC-Bayesian theory







< □ > < 同

(4) 国 (4) 日 (4) H (4) H

- PAC-Bayesian theory
 - Non-trivial generalization bound







→ < ∃→

3 / 23

- PAC-Bayesian theory
 - Non-trivial generalization bound
- Contributions







→ < ∃→

3 / 23

- PAC-Bayesian theory
 - Non-trivial generalization bound
- Contributions
 - Get rid of an approximation







▶ < ∃ ▶</p>

- PAC-Bayesian theory
 - Non-trivial generalization bound
- Contributions
 - Get rid of an approximation
 - Constant-time prediction relative to network depth







Plan of the presentation



▶ < Ξ >

Preliminary notions



< ∃→

- Preliminary notions
- e Neural networks aggregations



- Preliminary notions
- ② Neural networks aggregations
- Experimentation



- Preliminary notions
- Neural networks aggregations
- Experimentation
- Question period



Preliminary notions

Image: A test in te

< ∃ ▶

• Fully connected networks

- Fully connected networks
- Binary classification tasks

- Fully connected networks
- Binary classification tasks
- $F(\mathbf{x}) = (L_1 \circ \cdots \circ L_l)(\mathbf{x})$

- Fully connected networks
- Binary classification tasks
- $F(\mathbf{x}) = (L_1 \circ \cdots \circ L_l)(\mathbf{x})$
- $L_k(\mathbf{x}) = g_k(\mathbf{W}_k \mathbf{x} + \mathbf{b}_k) \ \forall k \in \{1, \dots, l\}$

- Fully connected networks
- Binary classification tasks
- $F(\mathbf{x}) = (L_1 \circ \cdots \circ L_l)(\mathbf{x})$

•
$$L_k(\mathbf{x}) = g_k(\mathbf{W}_k \mathbf{x} + \mathbf{b}_k) \ \forall k \in \{1, \dots, l\}$$





• Computation time savings



▶ ∢ ⊒

- Computation time savings
- Gains in memory required



- Computation time savings
- Gains in memory required
- Networks less prone to overfitting [1] [2]



J. Lin, C. Gan, and S. Han, "Defensive quantization: When efficiency meets robustness". CoRR, 2019
 M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1". NIPS, 2016.

Definition : The parts of a binary activated networks

Definition : The parts of a binary activated networks

• The first part is L_1

Definition : The parts of a binary activated networks

- The first part is L_1
- The second part is $(L_2 \circ \cdots \circ L_l)$

Binary activations - A visualization



Benjamin Leblanc (UL)

医水管医水管医 一座

Binary activations - A visualization



Benjamin Leblanc (UL)

医水管医水管医 一座

Binary activations - A visualization





Benjamin Leblanc (UL)

1 金属医金属医 画

Definition : The parts of a binary activated networks

- The first part (L₁)
- The second part $(L_2 \circ \cdots \circ L_l)$

We have $L_1: \mathcal{X} \to \{-1, +1\}^{d_1}$ and $(L_2 \circ \cdots \circ L_I): \{-1, +1\}^{d_1} \to \mathcal{Y}$

▶ < ∃ ▶</p>

Definition : The parts of a binary activated networks

- The first part (L_1) assigns a group to a given example
- The second part $(L_2 \circ \cdots \circ L_l)$

We have $L_1: \mathcal{X} \to \{-1, +1\}^{d_1}$ and $(L_2 \circ \cdots \circ L_I): \{-1, +1\}^{d_1} \to \mathcal{Y}$

Definition : The parts of a binary activated networks

- The first part (L_1) assigns a group to a given example
- The second part $(L_2 \circ \cdots \circ L_l)$ assigns a *prediction* to each group

We have $L_1: \mathcal{X} \to \{-1, +1\}^{d_1}$ and $(L_2 \circ \cdots \circ L_l): \{-1, +1\}^{d_1} \to \mathcal{Y}$

Neural networks aggregation

< ∃ ▶

Bayesian aspect

We suppose the weights $B = \langle \mathbf{W}_k \rangle_{k=1}^{l}$ to follow a Gaussian probability distribution, centred on B and with isotropic covariance matrix.

Bayesian aspect

We suppose the weights $B = \langle \mathbf{W}_k \rangle_{k=1}^l$ to follow a Gaussian probability distribution, centred on B and with isotropic covariance matrix.

The output of a neuron

The expected output of a sign function applied Gaussian distributed variable:

$$\underset{\textbf{v} \sim \mathcal{N}(\textbf{w},\textbf{I})}{\textbf{E}} \text{sgn}(\textbf{v} \cdot \textbf{a}) = \text{erf}\left(\frac{\textbf{w} \cdot \textbf{a}}{\sqrt{2}||\textbf{a}||}\right) \ ,$$

where $\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$ is the gaussian error function.

Bayesian aspect

We suppose the weights $B = \langle \mathbf{W}_k \rangle_{k=1}^l$ to follow a Gaussian probability distribution, centred on B and with isotropic covariance matrix.

The output of a neuron

The expected output of a sign function applied Gaussian distributed variable:

$$\underset{\mathbf{v} \sim \mathcal{N}(\mathbf{w},\mathbf{l})}{\mathsf{E}} \mathsf{sgn}(\mathbf{v} \cdot \mathbf{a}) = \mathsf{erf}\left(\frac{\mathbf{w} \cdot \mathbf{a}}{\sqrt{2}||\mathbf{a}||}\right) \;,$$

where $\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$ is the gaussian error function.



Probability of observing a given output for a given neuron

The binary activation acts as a Bernoulli distribution. $P(s_i \in \{-1, +1\})$ is given by

$$\Pr(\mathcal{L}_{k;i}(\mathbf{a}) = s) = \frac{1}{2} + \frac{s}{2} \operatorname{erf}\left(\frac{\mathbf{w}_{k,i} \cdot \mathbf{a}}{\sqrt{2}\|\mathbf{a}\|}\right), i \in \{1, \dots, d_k\}, k \in \{1, \dots, l\}.$$

Probability of observing a given output for a given neuron

The binary activation acts as a Bernoulli distribution. $P(s_i \in \{-1, +1\})$ is given by

$$\Pr(\mathcal{L}_{k;i}(\mathbf{a}) = s) = \frac{1}{2} + \frac{s}{2} \operatorname{erf}\left(\frac{\mathbf{w}_{k,i} \cdot \mathbf{a}}{\sqrt{2}\|\mathbf{a}\|}\right), i \in \{1, \dots, d_k\}, k \in \{1, \dots, l\}.$$

Probability of observing a given output for a given layer

Let $Pr((L_1 \circ \cdots \circ L_k)(\mathbf{a}) = \mathbf{s}) = L_{1,k}^{\mathbf{s}}(\mathbf{a})$. An iterative computation can be used:

$$\mathsf{Pr}(\mathcal{L}_{1,k}^{\mathbf{s}}(\mathbf{a})) = \begin{cases} \prod_{i=1}^{d_1} \mathsf{Pr}(\mathcal{L}_{1,i}^{s_i}(\mathbf{a})) & \text{si } k = 1, \\ \\ \sum_{\bar{\mathbf{s}}} \mathsf{Pr}(\mathcal{L}_{1,k}^{\bar{\mathbf{s}}}(\mathbf{a}) \mid \mathcal{L}_{1,k-1}^{\bar{\mathbf{s}}}(\mathbf{a})) \, \mathsf{Pr}(\mathcal{L}_{1,k-1}^{\bar{\mathbf{s}}}(\mathbf{a})) & \text{sinon.} \end{cases}$$

The exact computation of <u>ABNet</u>

Dynamic programming approach

The exact computation of <u>ABNet</u>

Dynamic programming approach

First layer:
$$\mathbf{P}_1(\mathbf{x}) = \left[\mathsf{Pr}(\mathcal{L}_1^{\mathsf{s}}(\mathbf{a})) \right]_{\mathbf{s} \in \mathcal{R}_1}$$

The exact computation of <u>ABNet</u>

Dynamic programming approach

First layer:
$$\mathbf{P}_{1}(\mathbf{x}) = \left[\Pr(\mathcal{L}_{1}^{s}(\mathbf{a})) \right]_{s \in R_{1}}$$

Subsequent layers: $\mathbf{P}_{k} = \left[\Pr(\mathcal{L}_{1,k-1}^{s}(\mathbf{a})) \right]_{s \in R_{k}}$
 $= \mathbf{\Psi}_{k} \cdot \mathbf{P}_{k-1}$,
with $\mathbf{\Psi}_{k} = \left[\Pr(\mathcal{L}_{1,k}^{s}(\mathbf{a}) \mid \mathcal{L}_{1,k-1}^{\overline{s}}(\mathbf{a})) \right]_{s \in R_{k}, \overline{s} \in R_{k-1}}$

Neural networks aggregation



(a) A typical neural network



(b) The architecture of **ABNet**, with underlying network from (a).

June 6th, 2023

Compact ABNet

$$\Pr(L_{1,l}^s(\mathbf{a})) = s \cdot (\Psi_l(\Psi_{l-1}(\ldots \Psi_3(\Psi_2 \mathbf{P}_1(\mathbf{x}))\ldots)))$$

Compact ABNet

$$\Pr(L_{1,l}^{s}(\mathbf{a})) = s \cdot (\Psi_{l}(\Psi_{l-1}(\dots \Psi_{3}(\Psi_{2}\mathbf{P}_{1}(\mathbf{x}))\dots)))$$
$$= \underbrace{s \cdot \Psi_{l} \cdot \Psi_{l-1} \cdot \dots \cdot \Psi_{3} \cdot \Psi_{2}}_{\mathbf{H}_{s}} \cdot \mathbf{P}_{1}(\mathbf{x}))$$

June 6th, 2023

16/23

Neural networks aggregation





(b) The architecture of **ABNet**, with underlying architecture from (a).



(c) The structure of **compact ABNet**, obtain from the network presented in (b).

• • = • • = •

Experimentations

▶ < ∃ ▶



(a) For different values of x, the probability of obtaining a certain representation as an output of the first hidden layer, obtained with $P_1(x)$, multiplied by H.



(a) For different values of x, the probability of obtaining a certain representation as an output of the first hidden layer, obtained with $P_1(x)$, multiplied by H.



(b) The output of the aggregation $Pr(L_{1,l})$.



(a) For different values of x, the probability of obtaining a certain representation as an output of the first hidden layer, obtained with $P_1(x)$, multiplied by H.



(b) The output of the aggregation $Pr(L_{1,l})$.



(c) ABNet as a classifier.



(a) For different values of x, the probability of obtaining a certain representation as an output of the first hidden layer, obtained with $P_1(x)$, multiplied by H.



(b) The output of the aggregation $Pr(L_{1,l})$.



(c) ABNet as a classifier.



(d) Unique network, centered on the *a posteriori* means of the Gaussian probability distribution.

CANAI 2023

19/23

ж



Figure: Impact of the depth for **PBGNet (dotted)** and **ABNet (continuous)** on the test error and the bound value depending on the width of the network on mnistLH. 5 random seeds.

Experimentations

Dataset	Model	L-1	d	Bound	Error _S	$Error_{\mathcal{T}}$
ads	PBGNet	3	2	0.192 ± 0.004	0.140 ± 0.004	0.141 ± 0.012
	ABNet	3	2	0.192 ± 0.004	0.140 ± 0.004	0.141 ± 0.012
	$PBGNet_\ell$	3	4	1.000 ± 0.001	0.018 ± 0.005	0.026 ± 0.004
	$ABNet_\ell$	3	4	0.887 ± 0.064	0.015 ± 0.003	0.026 ± 0.003
	EBP	2	2	_	0.003 ± 0.002	0.040 ± 0.008
	BC	1	4	_	0.025 ± 0.005	0.031 ± 0.004
	BNN	1	8	-	0.037 ± 0.002	0.038 ± 0.004
mnistLH	PBGNet	1	8	0.186 ± 0.028	0.091 ± 0.037	0.092 ± 0.036
	ABNet	3	4	0.162 ± 0.001	0.056 ± 0.001	0.058 ± 0.002
	$PBGNet_\ell$	3	8	1.000 ± 0.000	0.018 ± 0.003	0.038 ± 0.002
	$ABNet_\ell$	2	8	0.998 ± 0.003	0.025 ± 0.008	0.042 ± 0.006
	EBP	3	8	_	0.016 ± 0.002	0.043 ± 0.002
	BC	2	8	_	0.023 ± 0.002	0.035 ± 0.001
	BNN	1	2	_	0.123 ± 0.005	0.133 ± 0.004

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● □ ● ● ● ●

Questions / comments

> < 三 >

3

Thank you for listening :) !

Image: A test in te