

# SAE: Sequential Anchored Ensembles

---

Arnaud Delaunoy

December 9, 2021

University of Liège, Belgium

Advisor: Gilles Louppe



## Notations:

- $\mathbf{D}$ : Dataset
- $\theta$ : Neural network parameters
- $\mathbf{x}$ : Inputs
- $\mathbf{y}$ : Outputs

We want to compute  $p(\mathbf{y} | \mathbf{x}, \mathbf{D}) = \int p(\mathbf{y} | \mathbf{x}, \theta)p(\theta | \mathbf{D})d\theta$ ,

where  $p(\theta | \mathbf{D}) = \frac{p(\mathbf{D} | \theta)p(\theta)}{p(\mathbf{D})}$ .

**Training:** Ensemble of  $N$  neural networks such that  $\theta_{1,\dots,N}^* \sim p(\theta | \mathbf{D})$ .

**Prediction:**  $p(\mathbf{y} | \mathbf{x}, \mathbf{D}) \simeq \frac{1}{N} \sum_{i=1}^N p(\mathbf{y} | \mathbf{x}, \theta_i^*)$ .

**Idea:** Inject noise in the training procedure for the optima to be sampled from the Bayesian posterior

## Anchored Ensembling

**for**  $i$  in  $1, \dots, N$  **do**

$\theta_{\text{anc},i} \sim p(\theta)$  (Sample anchor)

$\theta_{\text{init},i} \leftarrow \text{init}()$  (Initialize NN)

$\theta_i^* \leftarrow \arg \max_{\theta} p(\mathbf{D} | \theta) p_{\text{anc},i}(\theta)$

**end for**

where  $p_{\text{anc},i} = \mathcal{N}(\theta_{\text{anc},i}, \Sigma_{\text{prior}})$ .

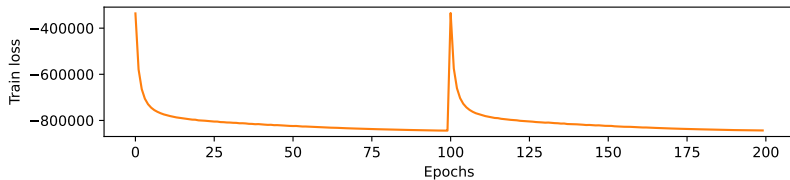
## Hypotheses:

- Normal prior:  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\text{prior}}, \boldsymbol{\Sigma}_{\text{prior}})$
- Normal likelihood  $p(\mathbf{D} | \boldsymbol{\theta})$  (also works for classification in practice)

If  $\boldsymbol{\theta}_{\text{anc}} \sim p(\boldsymbol{\theta})$  then  $\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} p(\mathbf{D} | \boldsymbol{\theta}) p_{\text{anc}}(\boldsymbol{\theta}) \sim p(\boldsymbol{\theta} | \mathbf{D})$   
(approximately).

# Anchored ensembles

Training an ensemble is computationally expensive.



# Sequential Anchored ensembles

If  $\theta_{\text{anc},i}$  is close to  $\theta_{\text{anc},i-1}$ , then  $\theta_i^*$  is close to  $\theta_{i-1}^*$

## Sequential Anchored Ensembling (SAE)

$\theta_{\text{anc},1} \sim p(\theta)$  (Sample first anchor)

$\theta_{\text{init},1} \leftarrow \text{init}()$  (Initialize NN)

$\theta_1^* \leftarrow \text{train}(\theta_{\text{anc},1}; \theta_{\text{init},1})$  (Long)

**for**  $i$  in  $2, \dots, M$  **do**

$\theta_{\text{anc},i} \leftarrow \text{mcmc\_step}(\theta_{\text{anc},i-1})$

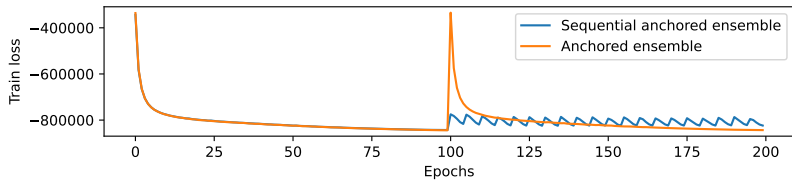
$\theta_{\text{init},i} \leftarrow \theta_{i-1}^*$

$\theta_i^* \leftarrow \text{train}(\theta_{\text{anc},i}; \theta_{\text{init},i})$  (Short)

**end for**

- Allow to build larger ensembles than AE
- SAE ensemble's members are correlated
- Can run SAE multiple times to benefit from different initializations

# Sequential Anchored ensembles





## Guided-walk Metropolis-Hastings

For SAE to work well, we need:

- $\theta_{\text{anc},i+1}$  close to  $\theta_{\text{anc},i}$  (short training)
- $\theta_{\text{anc},(1,\dots,M)}$  covers  $p(\theta)$  well

**Guided walk Metropolis-Hastings** [Gustafson, 1998]

## Guided walk Metropolis-Hastings

$$y \leftarrow \theta_{\text{anc},i-1} + d_{i-1}|z|, \quad z \sim \mathcal{N}(0, \sigma_{\text{step}})$$

$$\alpha \leftarrow \min\left(\frac{p(y)}{p(\theta_{\text{anc},i-1})}, 1\right)$$

$$u \sim \mathcal{U}(0, 1)$$

**if**  $u < \alpha$  **then**

$$\theta_{\text{anc},i} \leftarrow y$$

$$d_i \leftarrow d_{i-1}$$

**else**

$$\theta_{\text{anc},i} \leftarrow \theta_{\text{anc},i-1}$$

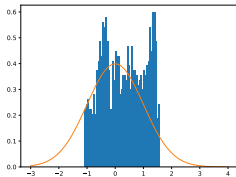
$$d_i \leftarrow -d_{i-1}$$

**end if**

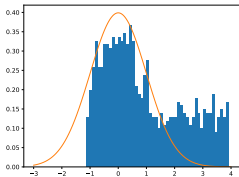
# Guided-walk Metropolis-Hastings

How to choose  $\sigma_{\text{step}}$  ?

(a)  $\sigma_{\text{step}} = 0.01$

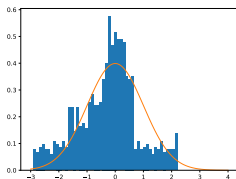


(b)  $\sigma_{\text{step}} = 0.02$

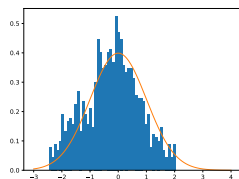


- Should be as small as possible
- Should span the prior  $\rightarrow$  we can verify this!

(c)  $\sigma_{\text{step}} = 0.03$

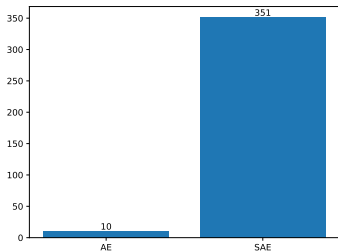


(d)  $\sigma_{\text{step}} = 0.05$

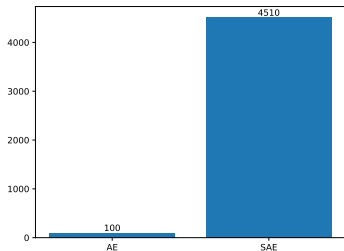


## Number of members in the ensemble

(a) 1000 epochs



(b) 10,000 epochs



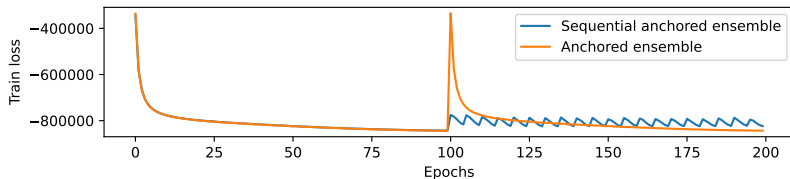
# Results

		Cifar10 Resnet		Cifar10-C Alexnet		IMDB		DermaMNIST		UCI-Gap
		Ag.	TV	Ag.	TV	Ag.	TV	Ag.	TV	$W_2$
1000 epochs	AE	0.849	0.201	0.726	0.262	<b>0.892</b>	<b>0.109</b>	0.877	0.104	<b>-0.148</b>
	SAE	<b>0.856</b>	<b>0.176</b>	<b>0.772</b>	<b>0.212</b>	0.887	0.110	<b>0.880</b>	<b>0.098</b>	-0.178
10,000 epochs	AE	0.862	0.199	0.746	0.236	<b>0.926</b>	<b>0.086</b>	<b>0.897</b>	0.089	<b>-0.137</b>
	SAE	<b>0.903</b>	<b>0.133</b>	<b>0.787</b>	<b>0.200</b>	0.916	0.099	0.893	<b>0.086</b>	-0.185

# Summary

**Mail:** a.delaunoy@uliege.be

**Twitter:** @ArnaudDelaunoy



- Tim Pearce, Felix Leibfried, and Alexandra Brintrup. Uncertainty in neural networks: Approximately bayesian ensembling. In International conference on artificial intelligence and statistics, pages 234–244. PMLR, 2020.
- Paul Gustafson. A guided walk metropolis algorithm. *Statistics and computing*, 8(4):357–364, 1998.