

Bayesian Compression for Deep Learning

Alexander
Markovich

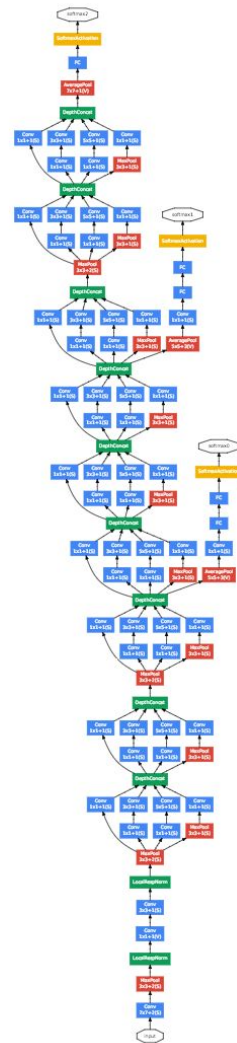


NATIONAL RESEARCH
UNIVERSITY

Moscow 2018

Motivation

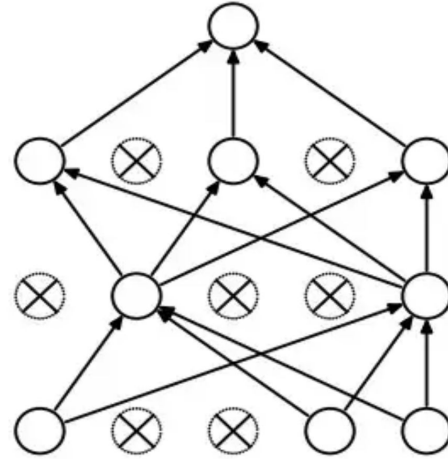
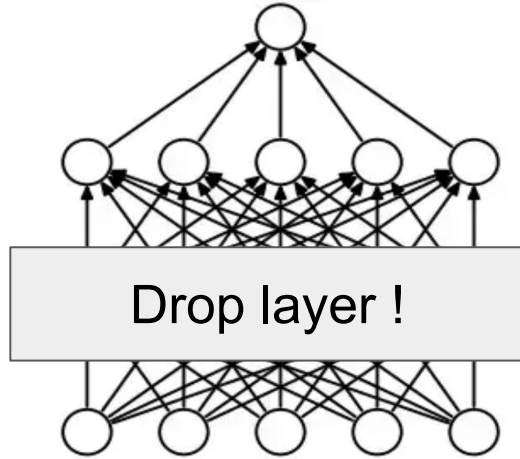
- Billions of parameters
- Slow inference
- Limited memory devices (smartphones, robots)
- Overfitting



Variety of methods

- Individual or structural sparsification / pruning
- Quantization
- Low rank approximation for weight matrices
- ...

Individual vs Structural



- Units
- Layers
- Filters
- ...

WE ARE THE BAYESIAN.

**YOU WILL BE ASSIMILATED. YOUR TECHNOLOGICAL DISTINCTIVENESS
WILL BE CONSIDERED A SPECIAL CASE OF OUR OWN. RESISTANCE IS FUTILE.**

Bayesian Inference <3



Given $\mathcal{D} = \left\{ (x_1, y_1), \dots, (x_n, y_n) \right\}$

Goal $p(y \mid x, w)$

and then?
$$p(w \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid w)p(w)}{p(\mathcal{D})}$$

Likelihood

Prior distribution

Evidence

(Stochastic) Variance Inference

$$p(w \mid \mathcal{D}) \approx q_\phi(w) \longleftarrow \text{Approximation}$$

$$\log p(\mathcal{D}) = \text{ELBO} + D_{KL}(q_\phi(w) \parallel p(w \mid \mathcal{D}))$$

$$\text{ELBO} = \underbrace{\sum_{i=1}^N \mathbb{E}_{q_\phi(w)} \log p(y_i \mid x_i, w)}_{\text{Data term}} - \underbrace{D_{KL}(q_\phi(w) \parallel p(w))}_{\text{Regularization}}$$

$$\min_{\phi} D_{KL}(q_\phi(w) \parallel p(w \mid \mathcal{D})) \Leftrightarrow \max_{\phi} \text{ELBO}$$

Relation to MDL

$$L(D) = \min_{H \in \mathcal{H}} \left\{ L(H) + L(D | H) \right\}$$

Hypothese

Data

$$\text{ELBO} = \underbrace{\mathbb{E}_{q_\phi(w)} \log p(\mathcal{D}|w)}_{\text{Error cost}} + \underbrace{\mathbb{E}_{q_\phi(w)} \log p(w) + \mathcal{H}(q_\phi(w))}_{\text{Complexity cost}}$$

Error cost

Complexity cost

(Sparse) Variational Dropout

$$B = (A \odot \Xi)W, \quad \xi_{mi} \sim p(\xi)$$
$$\begin{aligned} A &\in \mathbb{R}^{M \times I} \\ W &\in \mathbb{R}^{I \times O} \\ B &\in \mathbb{R}^{M \times O} \end{aligned}$$

$$\xi_{mi} \sim \text{Bernoulli}(1 - p) \longleftarrow \text{Binary Dropout}$$

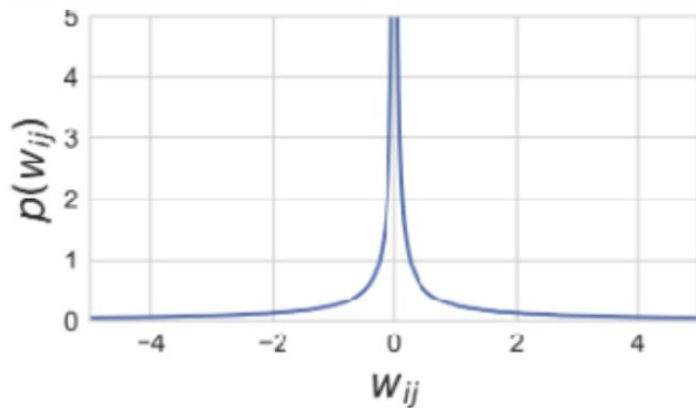
$$\xi_{mi} \sim \mathcal{N}(1, \alpha = \frac{p}{1-p}) \longleftarrow \text{Gaussian Dropout}$$

$$w_{ij} = \theta_{ij}\xi_{ij} = \theta_{ij}(1 + \sqrt{\alpha}\epsilon_{ij}) \sim \mathcal{N}(w_{ij} \mid \theta_{ij}, \alpha\theta_{ij}^2)$$
$$\epsilon_{ij} \sim \mathcal{N}(0, 1)$$

What about prior?



(Sparse) Variational Dropout

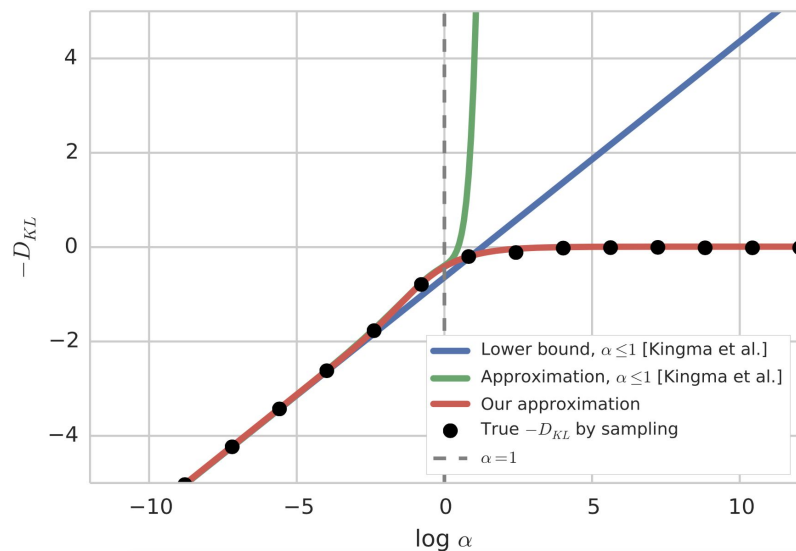


$$p(|w|) \propto \frac{1}{|w|}$$

$$q(w \mid \theta, \alpha) = \mathcal{N}(w \mid \theta, \alpha\theta^2)$$

(Sparse) Variational Dropout

$$\text{ELBO} = \sum_{i=1}^N \mathbb{E}_{q(w|\theta, \alpha)} \log p(y_i | x_i, w) - D_{KL}(q(w | \theta, \alpha) \| p(w))$$



$$-D_{KL}(q(w_{ij} | \theta_{ij}, \alpha_{ij}) \| p(w_{ij})) = \frac{1}{2} \log \alpha_{ij} - \mathbb{E}_{\epsilon \sim \mathcal{N}(1, \alpha_{ij})} \log |\epsilon| + C$$

(Sparse) Variational Dropout

$$\alpha_{ij} \rightarrow \infty$$

$$\theta_{ij} \rightarrow 0, \quad \alpha_{ij} \theta_{ij}^2 \rightarrow 0$$

$$\Downarrow$$

$$q(w_{ij} \mid \theta_{ij}, \alpha_{ij}) \rightarrow \mathcal{N}(w_{ij} \mid 0, 0) = \delta(w_{ij})$$

Hierarchical prior

Scale mixture of normals

$$w \sim \mathcal{N}(w \mid 0, z^2), \quad z \sim p(z)$$

Prior

$$p(w) = \int p(z) \mathcal{N}(w \mid 0, z^2) dz$$

$$p(z) \propto \frac{1}{|z|}$$

$$p(w) \propto \int \frac{1}{|z|} \mathcal{N}(w \mid 0, z^2) = \frac{1}{|w|}$$

$$p(z) \propto \sqrt{\det \mathcal{I}(z)}$$

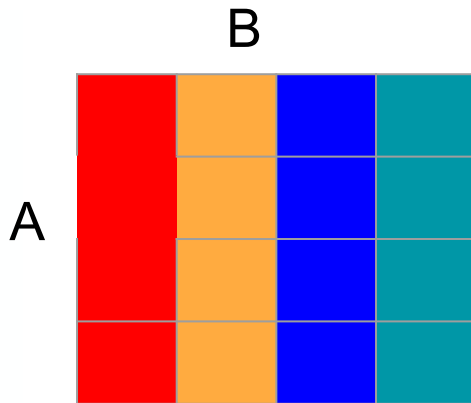
Group sparsity

$$p(W, z) \propto \prod_{i=1}^A \left[\frac{1}{|z_i|} \prod_j^B \mathcal{N}(w_{ij} \mid 0, z_i^2) \right]$$

$$q_\phi(W, z) = \prod_{i=1}^A \left[\mathcal{N}(z_i \mid \mu_{z_i}, \mu_{z_i}^2 \alpha_i) \prod_j^B \mathcal{N}(w_{ij} \mid z_i \mu_{ij}, z_i^2 \sigma_{ij}^2) \right]$$

Dropout rate per group

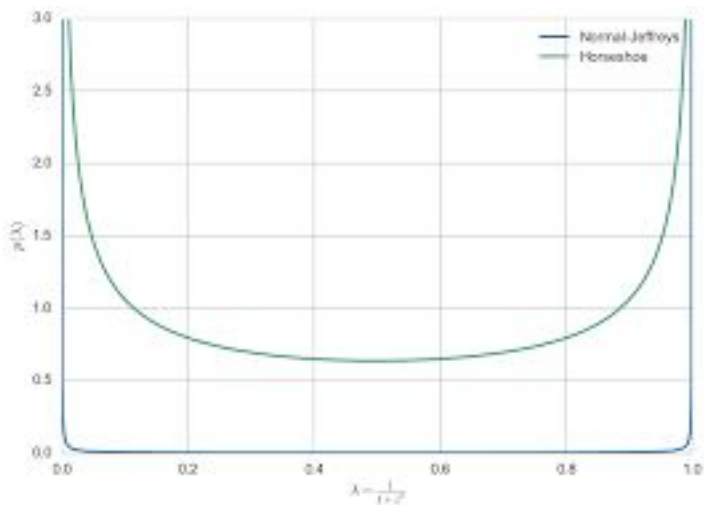
$\log \alpha_i > \text{threshold}$



Horseshoe prior

$$p(z) = \mathcal{C}^+(0, s) = 2 \left(s\pi(1 + (z/s)^2) \right)^{-1} \quad \text{half-Cauchy distribution}$$

$$s \sim \mathcal{C}^+(0, \tau_0); \quad \tilde{z}_i \sim \mathcal{C}^+(0, 1); \quad \tilde{w}_{ij} \sim \mathcal{N}(0, 1); \quad w_{ij} = \tilde{w}_{ij} \tilde{z}_i s,$$



Horseshoe prior

$$s \sim \mathcal{C}^+(0, \tau_0); \quad \tilde{z}_i \sim \mathcal{C}^+(0, 1); \quad \tilde{w}_{ij} \sim \mathcal{N}(0, 1); \quad w_{ij} = \tilde{w}_{ij} \tilde{z}_i s,$$

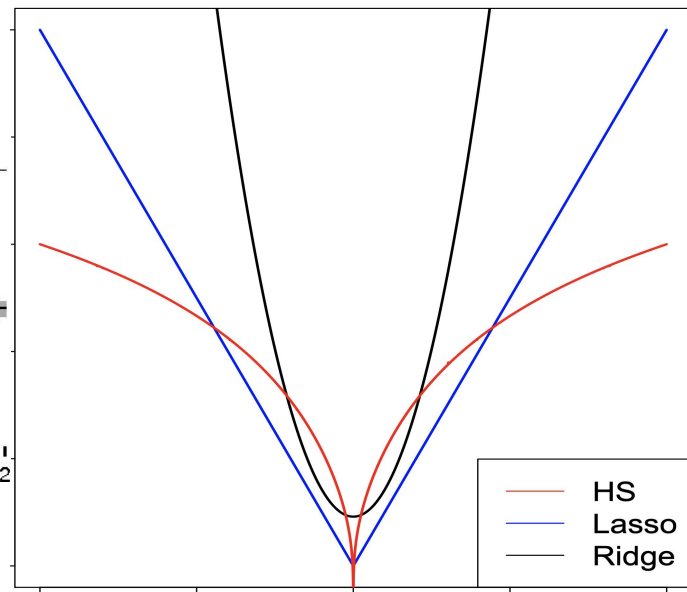
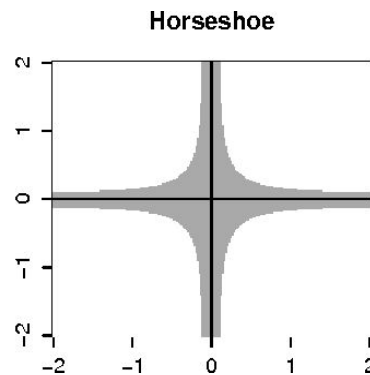
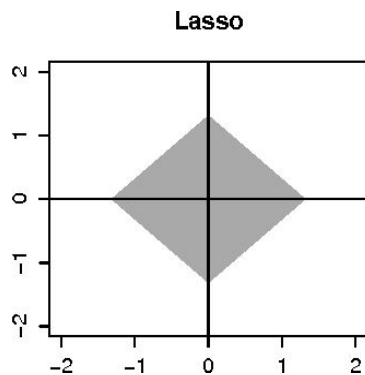
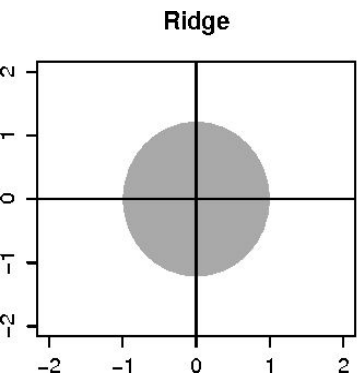


$$s_b \sim \mathcal{IG}(0.5, 1); \quad s_a \sim \mathcal{G}(0.5, \tau_0^2); \quad \tilde{\beta}_i \sim \mathcal{IG}(0.5, 1); \quad \tilde{\alpha}_i \sim \mathcal{G}(0.5, 1); \quad \tilde{w}_{ij} \sim \mathcal{N}(0, 1)$$

$$w_{ij} = \tilde{w}_{ij} \sqrt{s_a s_b \tilde{\alpha}_i \tilde{\beta}_i}.$$

Interesting fact

$\tilde{\alpha}_i, \tilde{\beta}_i \rightarrow 0 \Leftrightarrow$ Normal-Jeffreys prior = Horseshoe prior



Compression Rate

Model	Original Error %	Method	$\frac{ w \neq 0 }{ w } \%$	Compression Rates (Error %)	
				Fast Prediction	Maximum Compression
LeNet-300-100	1.6	DC	8.0	6 (1.6)	-
		DNS	1.8	28* (2.0)	-
		SWS	4.3	12* (1.9)	-
		Sparse VD	2.2	21(1.8)	84(1.8)
		BC-GNJ	10.8	9(1.8)	36(1.8)
		BC-GHS	10.6	9(1.8)	23(1.9)
LeNet-5-Caffe	0.9	DC	8.0	6*(0.7)	-
		DNS	0.9	55*(0.9)	-
		SWS	0.5	100*(1.0)	-
		Sparse VD	0.7	63(1.0)	228(1.0)
		BC-GNJ	0.9	108(1.0)	361(1.0)
		BC-GHS	0.6	156(1.0)	419(1.0)
VGG	8.4	BC-GNJ	6.7	14(8.6)	56(8.8)
		BC-GHS	5.5	18(9.0)	59(9.0)

Acceleration

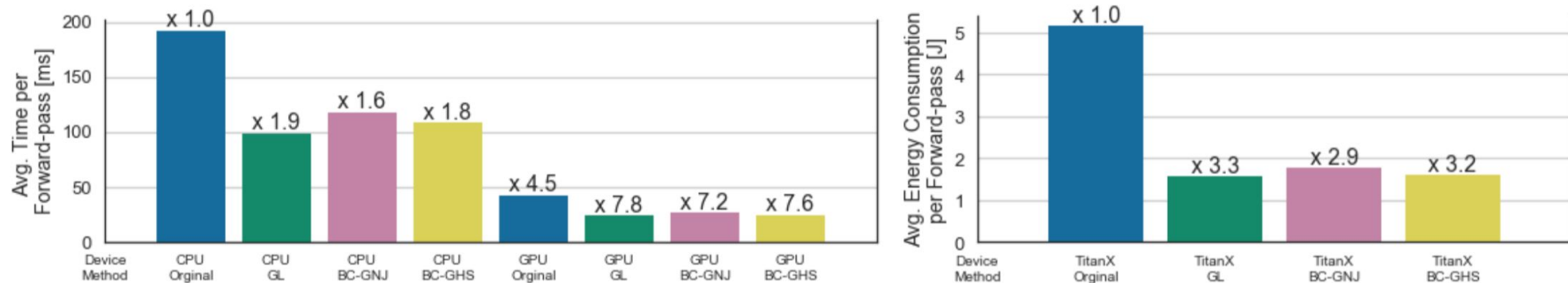


Figure 1: **Left:** Avg. Time a batch of 8192 samples takes to pass through LeNet-5-Caffe. Numbers on top of the bars represent speed-up factor relative to the CPU implementation of the original network. **Right:** Energy consumption of the GPU of the same process (when run on GPU).