

Deep Weight Prior

Andrei Atanov*, Arsenii Ashukha*,
Kirill Struminsky, Dmitry Vetrov, Max Welling

Bayesian Inference

Dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Likelihood $p(y_i | x_i, W)$

Prior $p(W)$

Bayes rule
$$p(W | \mathcal{D}) = \frac{p(\mathcal{D} | W)p(W)}{p(\mathcal{D})}$$

Variational Inference

Dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Likelihood $p(y_i | x_i, W)$

Prior $p(W)$

Variational posterior $q_\theta(W)$

$$D_{\text{KL}}(q_\theta(W) \| p(W | \mathcal{D})) \rightarrow \min_{\theta} \Leftrightarrow \mathcal{L}(\theta) \rightarrow \max_{\theta}$$

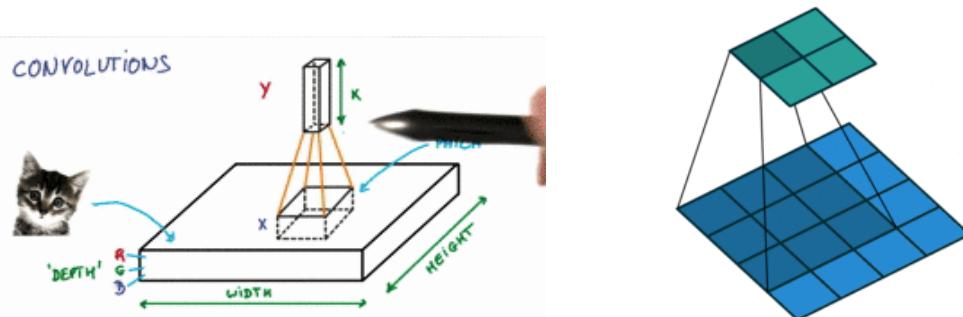
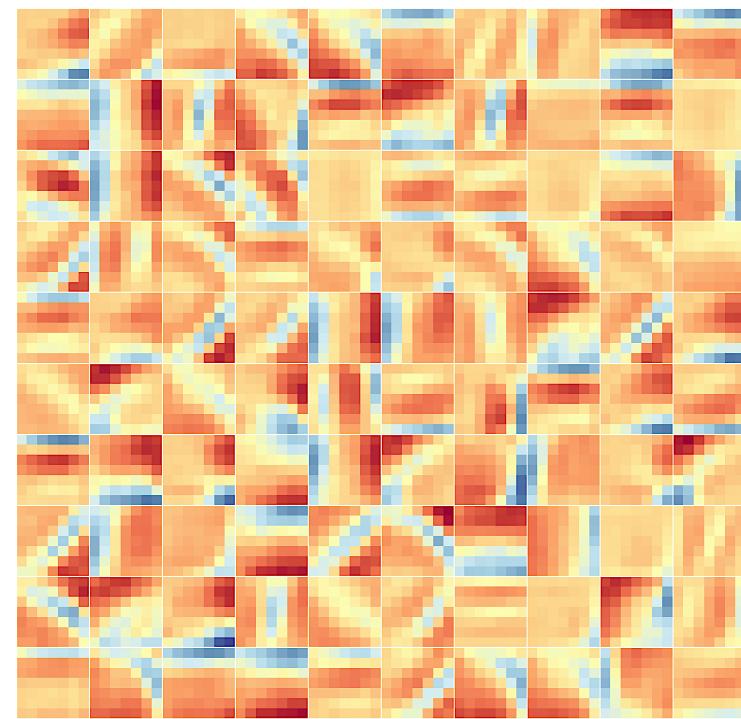
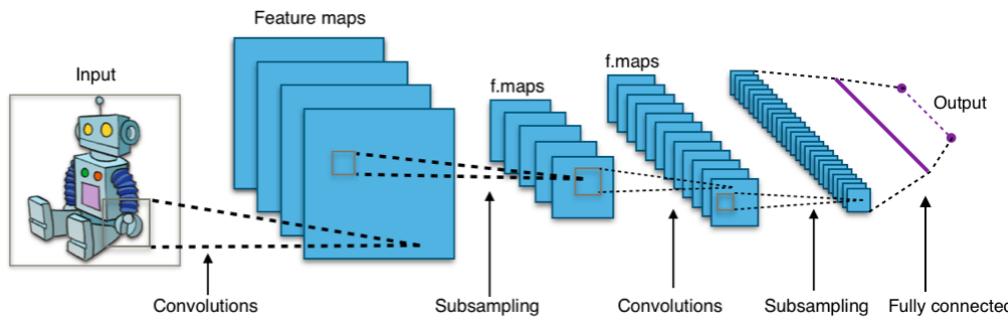
$$\text{ELBO} \quad \mathcal{L}(\theta) = \sum_{i=1}^N \mathbb{E}_{q_\theta(W)} \log p(y_i | x_i, W) - D_{\text{KL}}(q_\theta(W) \| p(W))$$

Prior distribution

Incorporate prior knowledge or specific properties into a model:

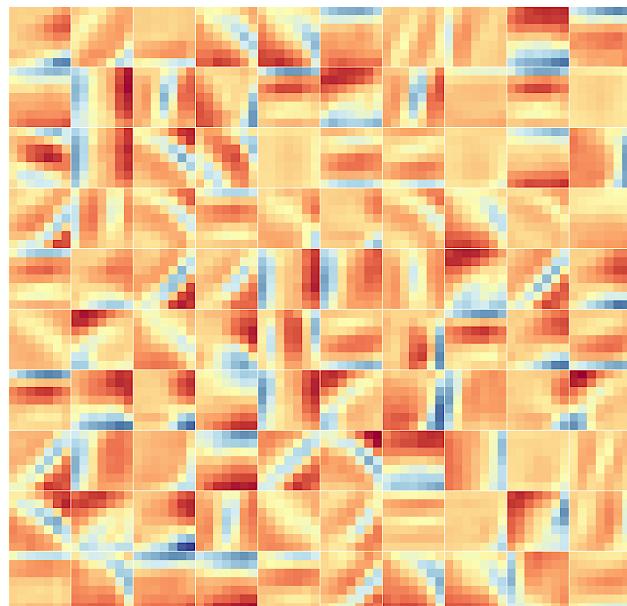
- Sparsification (Molchanov et al. 2017; Neklyudov et al. 2017)
 - Quantization (Ullrich et al., 2017)
 - Compression (Louizos et al., 2017)
-
- Usually fully-factorized
 - Does not take into account intrinsic structure of the weight

Convolutional Neural Networks

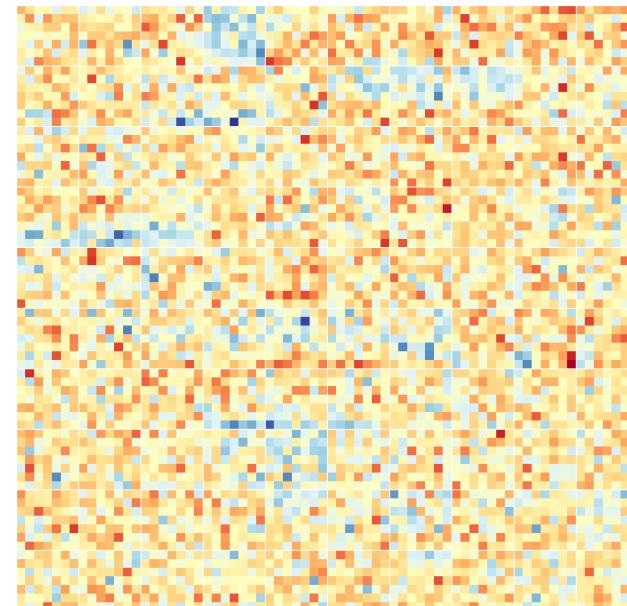


Learned convolutional kernels

Convolutional Neural Networks



Kernels learned on
large dataset



Kernels learned on
small dataset

Deep Weight Prior

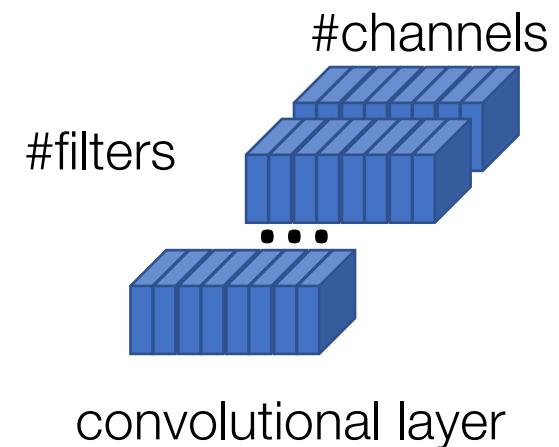
Consider prior factorized over layers, channels and filters:

$$p(W) = \prod_{l=1}^L \prod_{i=1}^{I_l} \prod_{j=1}^{O_l} p_l(w_{ij}^l)$$

#layers #channels #filters

L I_l O_l

w_{ij}^l kernel

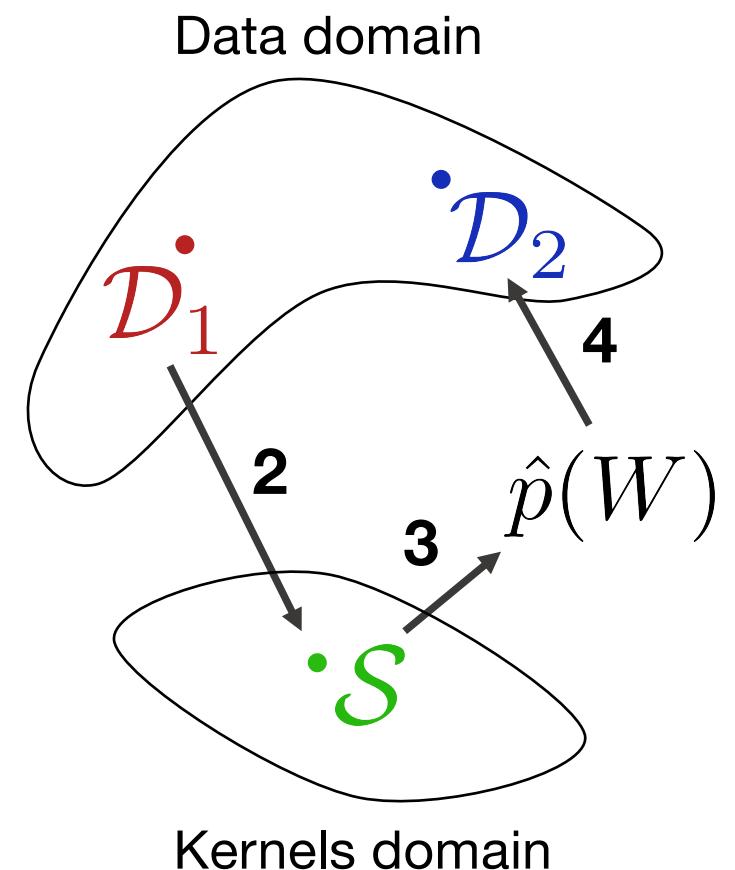


How to develop distribution that will **favor learned kernels?**

Generative models!

Learning Deep Weight Prior

1. Train CNNs on a given source dataset \mathcal{D}_1
2. Collect datasets \mathcal{S} of learned kernels
3. Train prior $\hat{p}(W)$ distributions on \mathcal{S}
4. Use this distribution for Bayesian inference on a new dataset \mathcal{D}_2



Generative models

Explicit models:

$$p(w)$$

- Kernel Density Estimation
- Normalizing Flows
- PixelCNN\RNN

- + Access to the density function
- High memory\computational cost

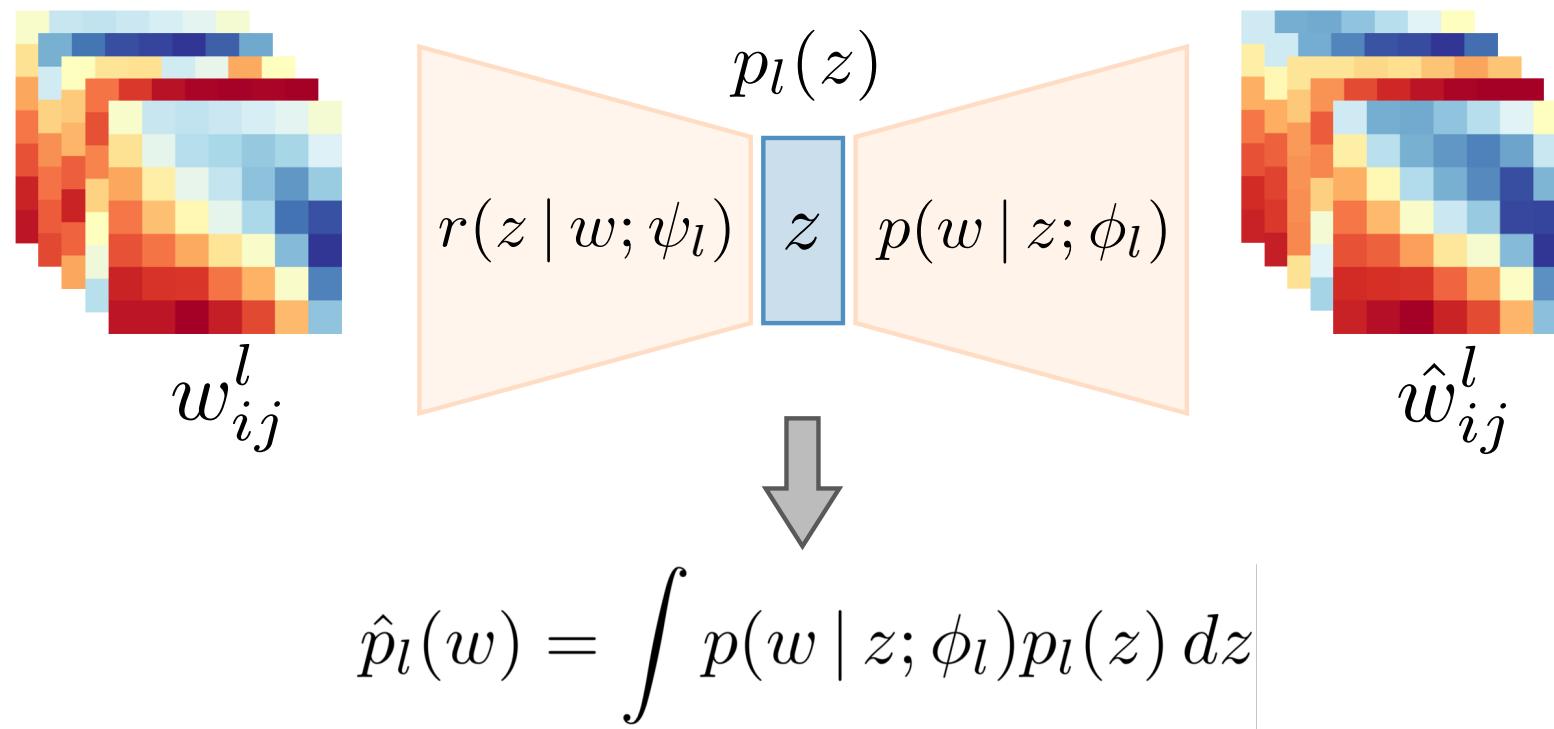
Implicit models:

$$p(w) = \int p(w \mid z)p(z)dz$$

- Variational auto-encoder

- + Memory and computational efficient
- No explicit form of the density

Learning Deep Weight Prior



Make ELBO tractable again

$$\mathcal{L}(\theta) = \sum_{i=1}^N \boxed{\mathbb{E}_{q_\theta(W)} \log p(y_i | x_i, W)} - \boxed{D_{\text{KL}}(q_\theta(W) \| p(W))}$$

Does not depend on a prior No more tractable for DWP!

KL-divergence for one kernel of l-th layer:

$$\begin{aligned} D_{\text{KL}}(q(w|\theta_{ij}^l) \| \hat{p}_l(w)) \\ = -H(q(w|\theta_{ij}^l)) - \mathbb{E}_{q(w|\theta_{ij}^l)} \boxed{\log \mathbb{E}_{p_l(z)} p(w|z; \phi_l)} \\ = \log \hat{p}_l(w) \end{aligned}$$

Make ELBO tractable again

$$D_{\text{KL}}(q(w) \parallel \hat{p}_l(w))$$

$$= -H(q(w)) - \mathbb{E}_{q(w)} \log \hat{p}_l(w)$$

$$= -H(q(w)) - \mathbb{E}_{q(w)} \mathbb{E}_{r(z|w)} \log \frac{p(w \mid z; \phi_l) p_l(z)}{p(z|w)} \frac{r(z \mid w)}{r(z \mid w)}$$

$$= -H(q(w)) + \mathbb{E}_{q(w)} \left[D_{\text{KL}}(r(z|w) \parallel p_l(z)) - \mathbb{E}_{r(z|w)} \log p(w|z; \phi_l) \right]$$

$$- \mathbb{E}_{q(w)} D_{\text{KL}}(r(z|w) \parallel p(z|w))$$

Make ELBO tractable again

$$D_{\text{KL}}(q(w) \| \hat{p}_l(w))$$

$$= -H(q(w)) - \mathbb{E}_{q(w)} \log \hat{p}_l(w)$$

$$= -H(q(w)) - \mathbb{E}_{q(w)} \mathbb{E}_{r(z|w)} \log \frac{p(w | z; \phi_l) p_l(z)}{p(z|w)} \frac{r(z | w)}{r(z | w)}$$

$$= \boxed{-H(q(w)) + \mathbb{E}_{q(w)} [D_{\text{KL}}(r(z|w) \| p_l(z)) - \mathbb{E}_{r(z|w)} \log p(w|z; \phi_l)]}$$

KL upper bound

$$\boxed{-\mathbb{E}_{q(w)} D_{\text{KL}}(r(z|w) \| p(z|w))}$$

<= 0

Make ELBO tractable again

- **Upper bound on the whole KL-term:**

$$D_{\text{KL}}(q(W) \| \hat{p}(W)) = \sum_{l,i,j} D_{\text{KL}}(q(w_{ij}^l | \theta_{ij}^l) \| \hat{p}_l(w_{ij}^l)) \leq \sum_{l,i,j} (-H(q(w_{ij}^l | \theta_{ij}^l)) + \\ + \mathbb{E}_{q(w_{ij}^l | \theta_{ij}^l)} [D_{\text{KL}}(r(z | w_{ij}^l; \psi_l) \| p_l(z)) - \mathbb{E}_{r(z | w_{ij}^l; \psi_l)} \log p(w_{ij}^l | z; \phi_l)]) = D_{\text{KL}}^{\text{bound}}$$

- **Final auxiliary lower bound:**

$$\mathcal{L}^{\text{aux}}(\theta, \psi) = L_D - D_{\text{KL}}^{\text{bound}} \leq L_D - D_{\text{KL}}(q_\theta(W) \| \hat{p}(W)) = \mathcal{L}(\theta)$$

$$\mathcal{L}(\theta, \psi) \rightarrow \max_{\theta, \psi}$$

Make ELBO tractable again

- **Upper bound on the whole KL-term:**

$$D_{\text{KL}}(q(W) \| \hat{p}(W)) = \sum_{l,i,j} D_{\text{KL}}(q(w_{ij}^l | \theta_{ij}^l) \| \hat{p}_l(w_{ij}^l)) \leq \sum_{l,i,j} (-H(q(w_{ij}^l | \theta_{ij}^l)) + \\ + \mathbb{E}_{q(w_{ij}^l | \theta_{ij}^l)} [D_{\text{KL}}(r(z | w_{ij}^l; \psi_l) \| p_l(z)) - \mathbb{E}_{r(z | w_{ij}^l; \psi_l)} \log p(w_{ij}^l | z; \phi_l)]) = D_{\text{KL}}^{\text{bound}}$$

- **Final auxiliary lower bound:**

$$\mathcal{L}^{\text{aux}}(\theta, \psi) = L_D - D_{\text{KL}}^{\text{bound}} \leq L_D - D_{\text{KL}}(q_\theta(W) \| \hat{p}(W)) = \mathcal{L}(\theta)$$

$$\mathcal{L}(\theta, \psi) \rightarrow \max_{\theta, \psi}$$

ϕ **is fixed!**

Algorithm 1 Stochastic Variational Inference With Implicit Prior Distribution

Require: dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Require: variational approximations $q(w | \theta_{ij}^l)$ and reverse models $r(z | w; \psi_l)$

Require: reconstruction models $p(w | z; \phi_l)$, priors for auxiliary variables $p_l(z)$

while not converged **do**

$\hat{M} \leftarrow$ mini-batch of objects form dataset \mathcal{D}

$\hat{w}_{ij}^l \leftarrow$ sample weights from $q(w | \theta_{ij}^l)$ with reparametrization

$\hat{z}_{ij}^l \leftarrow$ sample auxiliary variables from $r(z | \hat{w}_{ij}^l; \psi_l)$ with reparametrization

$\hat{\mathcal{L}}^{aux} \leftarrow L_{\hat{M}} + \sum_{l,i,j} -\log q(\hat{w}_{ij}^l | \theta_{ij}^l) - \log r(\hat{z}_{ij}^l | \hat{w}_{ij}^l; \psi_l) + \log p_l(\hat{z}_{ij}^l) + \log p(\hat{w}_{ij}^l | \hat{z}_{ij}^l; \phi_l)$

 Obtain unbiased estimate \hat{g} with $\mathbb{E}[\hat{g}] = \nabla \hat{\mathcal{L}}^{aux}$ by differentiating $\hat{\mathcal{L}}^{aux}$

 Update parameters θ and ψ using gradient \hat{g} and a stochastic optimization algorithm

end while

return Parameters θ, ψ

Algorithm 1 Stochastic Variational Inference With Implicit Prior Distribution

Require: dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ 

Require: variational approximations $q(w | \theta_{ij}^l)$ and reverse models $r(z | w; \psi_l)$

Require: reconstruction models $p(w | z; \phi_l)$, priors for auxiliary variables $p_l(z)$

while not converged **do**

$\hat{M} \leftarrow$ mini-batch of objects form dataset \mathcal{D}

$\hat{w}_{ij}^l \leftarrow$ sample weights from $q(w | \theta_{ij}^l)$ with reparametrization

$\hat{z}_{ij}^l \leftarrow$ sample auxiliary variables from $r(z | \hat{w}_{ij}^l; \psi_l)$ with reparametrization

$\hat{\mathcal{L}}^{aux} \leftarrow L_{\hat{M}} + \sum_{l,i,j} -\log q(\hat{w}_{ij}^l | \theta_{ij}^l) - \log r(\hat{z}_{ij}^l | \hat{w}_{ij}^l; \psi_l) + \log p_l(\hat{z}_{ij}^l) + \log p(\hat{w}_{ij}^l | \hat{z}_{ij}^l; \phi_l)$

 Obtain unbiased estimate \hat{g} with $\mathbb{E}[\hat{g}] = \nabla \hat{\mathcal{L}}^{aux}$ by differentiating $\hat{\mathcal{L}}^{aux}$

 Update parameters θ and ψ using gradient \hat{g} and a stochastic optimization algorithm

end while

return Parameters θ, ψ

Algorithm 1 Stochastic Variational Inference With Implicit Prior Distribution

Require: dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Require: variational approximations $q(w | \theta_{ij}^l)$ and reverse models $r(z | w; \psi_l)$ 

Require: reconstruction models $p(w | z; \phi_l)$, priors for auxiliary variables $p_l(z)$

while not converged **do**

$\hat{M} \leftarrow$ mini-batch of objects form dataset \mathcal{D}

$\hat{w}_{ij}^l \leftarrow$ sample weights from $q(w | \theta_{ij}^l)$ with reparametrization

$\hat{z}_{ij}^l \leftarrow$ sample auxiliary variables from $r(z | \hat{w}_{ij}^l; \psi_l)$ with reparametrization

$\hat{\mathcal{L}}^{aux} \leftarrow L_{\hat{M}} + \sum_{l,i,j} -\log q(\hat{w}_{ij}^l | \theta_{ij}^l) - \log r(\hat{z}_{ij}^l | \hat{w}_{ij}^l; \psi_l) + \log p_l(\hat{z}_{ij}^l) + \log p(\hat{w}_{ij}^l | \hat{z}_{ij}^l; \phi_l)$

 Obtain unbiased estimate \hat{g} with $\mathbb{E}[\hat{g}] = \nabla \hat{\mathcal{L}}^{aux}$ by differentiating $\hat{\mathcal{L}}^{aux}$

 Update parameters θ and ψ using gradient \hat{g} and a stochastic optimization algorithm

end while

return Parameters θ, ψ

Algorithm 1 Stochastic Variational Inference With Implicit Prior Distribution

Require: dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Require: variational approximations $q(w | \theta_{ij}^l)$ and reverse models $r(z | w; \psi_l)$

Require: reconstruction models $p(w | z; \phi_l)$, priors for auxiliary variables $p_l(z)$ 

while not converged **do**

$\hat{M} \leftarrow$ mini-batch of objects form dataset \mathcal{D}

$\hat{w}_{ij}^l \leftarrow$ sample weights from $q(w | \theta_{ij}^l)$ with reparametrization

$\hat{z}_{ij}^l \leftarrow$ sample auxiliary variables from $r(z | \hat{w}_{ij}^l; \psi_l)$ with reparametrization

$\hat{\mathcal{L}}^{aux} \leftarrow L_{\hat{M}} + \sum_{l,i,j} -\log q(\hat{w}_{ij}^l | \theta_{ij}^l) - \log r(\hat{z}_{ij}^l | \hat{w}_{ij}^l; \psi_l) + \log p_l(\hat{z}_{ij}^l) + \log p(\hat{w}_{ij}^l | \hat{z}_{ij}^l; \phi_l)$

 Obtain unbiased estimate \hat{g} with $\mathbb{E}[\hat{g}] = \nabla \hat{\mathcal{L}}^{aux}$ by differentiating $\hat{\mathcal{L}}^{aux}$

 Update parameters θ and ψ using gradient \hat{g} and a stochastic optimization algorithm

end while

return Parameters θ, ψ

Algorithm 1 Stochastic Variational Inference With Implicit Prior Distribution

Require: dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Require: variational approximations $q(w | \theta_{ij}^l)$ and reverse models $r(z | w; \psi_l)$

Require: reconstruction models $p(w | z; \phi_l)$, priors for auxiliary variables $p_l(z)$

while not converged **do**

 → $\hat{M} \leftarrow$ mini-batch of objects form dataset \mathcal{D}

$\hat{w}_{ij}^l \leftarrow$ sample weights from $q(w | \theta_{ij}^l)$ with reparametrization

$\hat{z}_{ij}^l \leftarrow$ sample auxiliary variables from $r(z | \hat{w}_{ij}^l; \psi_l)$ with reparametrization

$\hat{\mathcal{L}}^{aux} \leftarrow L_{\hat{M}} + \sum_{l,i,j} -\log q(\hat{w}_{ij}^l | \theta_{ij}^l) - \log r(\hat{z}_{ij}^l | \hat{w}_{ij}^l; \psi_l) + \log p_l(\hat{z}_{ij}^l) + \log p(\hat{w}_{ij}^l | \hat{z}_{ij}^l; \phi_l)$

 Obtain unbiased estimate \hat{g} with $\mathbb{E}[\hat{g}] = \nabla \hat{\mathcal{L}}^{aux}$ by differentiating $\hat{\mathcal{L}}^{aux}$

 Update parameters θ and ψ using gradient \hat{g} and a stochastic optimization algorithm

end while

return Parameters θ, ψ

Algorithm 1 Stochastic Variational Inference With Implicit Prior Distribution

Require: dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Require: variational approximations $q(w | \theta_{ij}^l)$ and reverse models $r(z | w; \psi_l)$

Require: reconstruction models $p(w | z; \phi_l)$, priors for auxiliary variables $p_l(z)$

while not converged **do**

$\hat{M} \leftarrow$ mini-batch of objects form dataset \mathcal{D}

 → $\hat{w}_{ij}^l \leftarrow$ sample weights from $q(w | \theta_{ij}^l)$ with reparametrization

$\hat{z}_{ij}^l \leftarrow$ sample auxiliary variables from $r(z | \hat{w}_{ij}^l; \psi_l)$ with reparametrization

$\hat{\mathcal{L}}^{aux} \leftarrow L_{\hat{M}} + \sum_{l,i,j} -\log q(\hat{w}_{ij}^l | \theta_{ij}^l) - \log r(\hat{z}_{ij}^l | \hat{w}_{ij}^l; \psi_l) + \log p_l(\hat{z}_{ij}^l) + \log p(\hat{w}_{ij}^l | \hat{z}_{ij}^l; \phi_l)$

 Obtain unbiased estimate \hat{g} with $\mathbb{E}[\hat{g}] = \nabla \hat{\mathcal{L}}^{aux}$ by differentiating $\hat{\mathcal{L}}^{aux}$

 Update parameters θ and ψ using gradient \hat{g} and a stochastic optimization algorithm

end while

return Parameters θ, ψ

Algorithm 1 Stochastic Variational Inference With Implicit Prior Distribution

Require: dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Require: variational approximations $q(w | \theta_{ij}^l)$ and reverse models $r(z | w; \psi_l)$

Require: reconstruction models $p(w | z; \phi_l)$, priors for auxiliary variables $p_l(z)$

while not converged **do**

$\hat{M} \leftarrow$ mini-batch of objects form dataset \mathcal{D}

$\hat{w}_{ij}^l \leftarrow$ sample weights from $q(w | \theta_{ij}^l)$ with reparametrization

 → $\hat{z}_{ij}^l \leftarrow$ sample auxiliary variables from $r(z | \hat{w}_{ij}^l; \psi_l)$ with reparametrization

$\hat{\mathcal{L}}^{aux} \leftarrow L_{\hat{M}} + \sum_{l,i,j} -\log q(\hat{w}_{ij}^l | \theta_{ij}^l) - \log r(\hat{z}_{ij}^l | \hat{w}_{ij}^l; \psi_l) + \log p_l(\hat{z}_{ij}^l) + \log p(\hat{w}_{ij}^l | \hat{z}_{ij}^l; \phi_l)$

 Obtain unbiased estimate \hat{g} with $\mathbb{E}[\hat{g}] = \nabla \hat{\mathcal{L}}^{aux}$ by differentiating $\hat{\mathcal{L}}^{aux}$

 Update parameters θ and ψ using gradient \hat{g} and a stochastic optimization algorithm

end while

return Parameters θ, ψ

Algorithm 1 Stochastic Variational Inference With Implicit Prior Distribution

Require: dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Require: variational approximations $q(w | \theta_{ij}^l)$ and reverse models $r(z | w; \psi_l)$

Require: reconstruction models $p(w | z; \phi_l)$, priors for auxiliary variables $p_l(z)$

while not converged **do**

$\hat{M} \leftarrow$ mini-batch of objects form dataset \mathcal{D}

$\hat{w}_{ij}^l \leftarrow$ sample weights from $q(w | \theta_{ij}^l)$ with reparametrization

$\hat{z}_{ij}^l \leftarrow$ sample auxiliary variables from $r(z | \hat{w}_{ij}^l; \psi_l)$ with reparametrization

$\hat{\mathcal{L}}^{aux} \leftarrow L_{\hat{M}} + \sum_{l,i,j} -\log q(\hat{w}_{ij}^l | \theta_{ij}^l) - \log r(\hat{z}_{ij}^l | \hat{w}_{ij}^l; \psi_l) + \log p_l(\hat{z}_{ij}^l) + \log p(\hat{w}_{ij}^l | \hat{z}_{ij}^l; \phi_l)$

 Obtain unbiased estimate \hat{g} with $\mathbb{E}[\hat{g}] = \nabla \hat{\mathcal{L}}^{aux}$ by differentiating $\hat{\mathcal{L}}^{aux}$

 Update parameters θ and ψ using gradient \hat{g} and a stochastic optimization algorithm

end while

return Parameters θ, ψ

Algorithm 1 Stochastic Variational Inference With Implicit Prior Distribution

Require: dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Require: variational approximations $q(w | \theta_{ij}^l)$ and reverse models $r(z | w; \psi_l)$

Require: reconstruction models $p(w | z; \phi_l)$, priors for auxiliary variables $p_l(z)$

while not converged **do**

$\hat{M} \leftarrow$ mini-batch of objects form dataset \mathcal{D}

$\hat{w}_{ij}^l \leftarrow$ sample weights from $q(w | \theta_{ij}^l)$ with reparametrization

$\hat{z}_{ij}^l \leftarrow$ sample auxiliary variables from $r(z | \hat{w}_{ij}^l; \psi_l)$ with reparametrization

$\hat{\mathcal{L}}^{aux} \leftarrow L_{\hat{M}} + \sum_{l,i,j} -\log q(\hat{w}_{ij}^l | \theta_{ij}^l) - \log r(\hat{z}_{ij}^l | \hat{w}_{ij}^l; \psi_l) + \log p_l(\hat{z}_{ij}^l) + \log p(\hat{w}_{ij}^l | \hat{z}_{ij}^l; \phi_l)$

 Obtain unbiased estimate \hat{g} with $\mathbb{E}[\hat{g}] = \nabla \hat{\mathcal{L}}^{aux}$ by differentiating $\hat{\mathcal{L}}^{aux}$

 Update parameters θ and ψ using gradient \hat{g} and a stochastic optimization algorithm

end while

return Parameters θ, ψ

$$-H(q(w)) + \mathbb{E}_{q(w)} [D_{KL}(r(z|w) \| p_l(z)) - \mathbb{E}_{r(z|w)} \log p(w|z; \phi_l)]$$

KL upper bound

Algorithm 1 Stochastic Variational Inference With Implicit Prior Distribution

Require: dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Require: variational approximations $q(w | \theta_{ij}^l)$ and reverse models $r(z | w; \psi_l)$

Require: reconstruction models $p(w | z; \phi_l)$, priors for auxiliary variables $p_l(z)$

while not converged **do**

$\hat{M} \leftarrow$ mini-batch of objects form dataset \mathcal{D}

$\hat{w}_{ij}^l \leftarrow$ sample weights from $q(w | \theta_{ij}^l)$ with reparametrization

$\hat{z}_{ij}^l \leftarrow$ sample auxiliary variables from $r(z | \hat{w}_{ij}^l; \psi_l)$ with reparametrization

$\hat{\mathcal{L}}^{aux} \leftarrow L_{\hat{M}} + \sum_{l,i,j} -\log q(\hat{w}_{ij}^l | \theta_{ij}^l) - \log r(\hat{z}_{ij}^l | \hat{w}_{ij}^l; \psi_l) + \log p_l(\hat{z}_{ij}^l) + \log p(\hat{w}_{ij}^l | \hat{z}_{ij}^l; \phi_l)$

 Obtain unbiased estimate \hat{g} with $\mathbb{E}[\hat{g}] = \nabla \hat{\mathcal{L}}^{aux}$ by differentiating $\hat{\mathcal{L}}^{aux}$

 Update parameters θ and ψ using gradient \hat{g} and a stochastic optimization algorithm

end while

return Parameters θ, ψ

$$-H(q(w)) + \mathbb{E}_{q(w)} [D_{KL}(r(z|w) \| p_l(z)) - \mathbb{E}_{r(z|w)} \log p(w|z; \phi_l)]$$

KL upper bound

Algorithm 1 Stochastic Variational Inference With Implicit Prior Distribution

Require: dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Require: variational approximations $q(w | \theta_{ij}^l)$ and reverse models $r(z | w; \psi_l)$

Require: reconstruction models $p(w | z; \phi_l)$, priors for auxiliary variables $p_l(z)$

while not converged **do**

$\hat{M} \leftarrow$ mini-batch of objects form dataset \mathcal{D}

$\hat{w}_{ij}^l \leftarrow$ sample weights from $q(w | \theta_{ij}^l)$ with reparametrization

$\hat{z}_{ij}^l \leftarrow$ sample auxiliary variables from $r(z | \hat{w}_{ij}^l; \psi_l)$ with reparametrization

$\hat{\mathcal{L}}^{aux} \leftarrow L_{\hat{M}} + \sum_{l,i,j} -\log q(\hat{w}_{ij}^l | \theta_{ij}^l) - \log r(\hat{z}_{ij}^l | \hat{w}_{ij}^l; \psi_l) + \log p_l(\hat{z}_{ij}^l) + \log p(\hat{w}_{ij}^l | \hat{z}_{ij}^l; \phi_l)$

 Obtain unbiased estimate \hat{g} with $\mathbb{E}[\hat{g}] = \nabla \hat{\mathcal{L}}^{aux}$ by differentiating $\hat{\mathcal{L}}^{aux}$

 Update parameters θ and ψ using gradient \hat{g} and a stochastic optimization algorithm

end while

return Parameters θ, ψ

$$-H(q(w)) + \mathbb{E}_{q(w)} [D_{KL}(r(z|w) \| p_l(z)) - \mathbb{E}_{r(z|w)} \log p(w|z; \phi_l)]$$

KL upper bound

Experiments

Train prior on:



\mathcal{D}_1 notMNIST

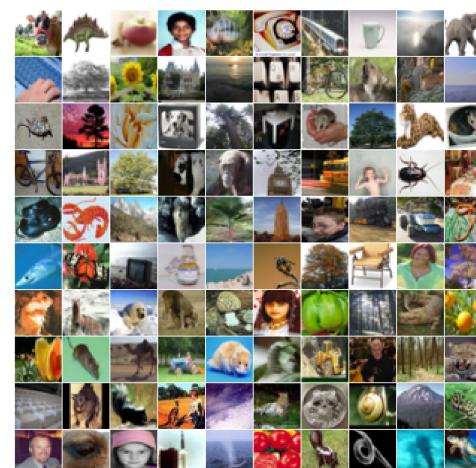
$$\hat{p}(W)$$

Use prior for:

5	0	4	1	9	2	1	3	1	4
3	5	3	6	1	7	2	8	6	9
4	0	9	1	1	2	4	3	2	7
3	8	6	9	0	5	6	0	7	6
1	8	7	1	9	3	9	8	5	9
3	0	7	4	9	8	0	9	4	1
4	4	6	0	4	5	6	1	0	0
1	7	1	6	3	0	2	1	1	7
9	0	2	6	7	8	3	9	0	4
6	7	4	6	8	0	7	8	3	1

\mathcal{D}_2 small
MNIST

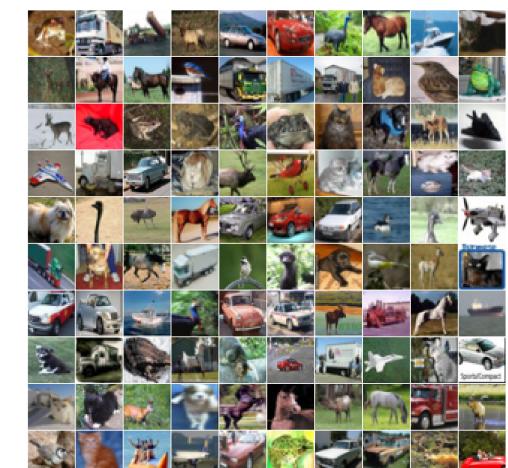
Train prior on:



\mathcal{D}_1 CIFAR100

$$\hat{p}(W)$$

Use prior for:

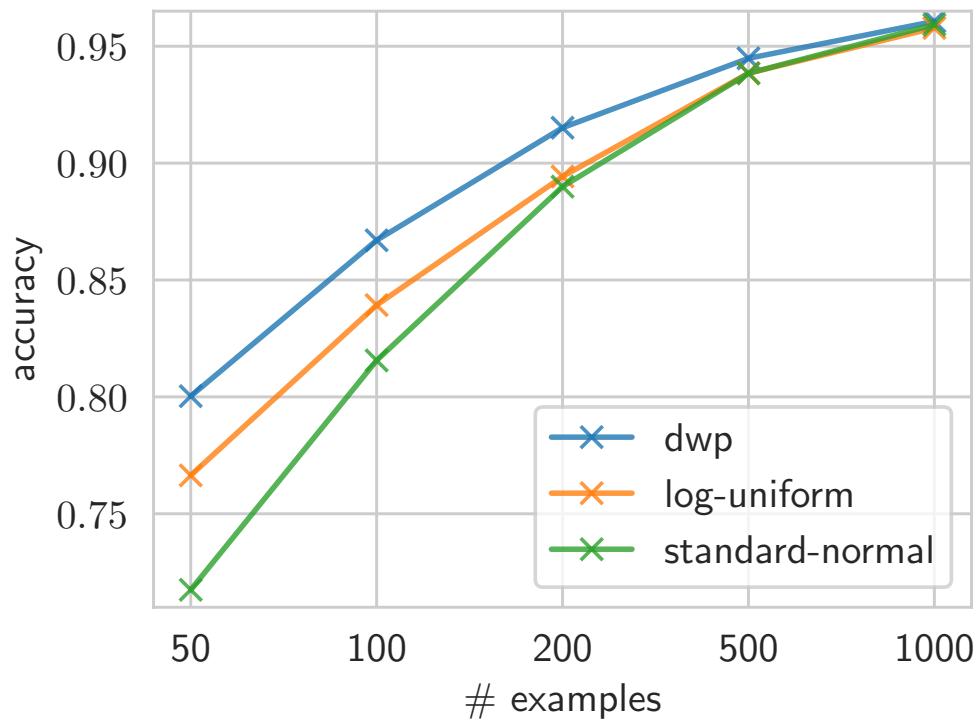


\mathcal{D}_2 small
CIFAR10

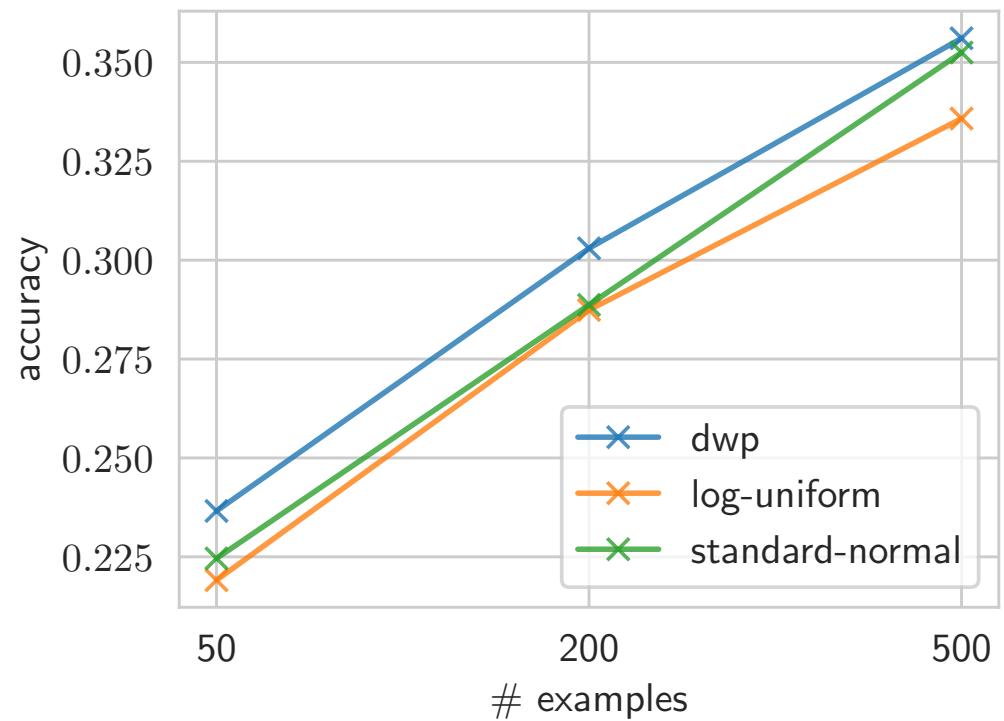
$$\hat{p}(W)$$

Classification

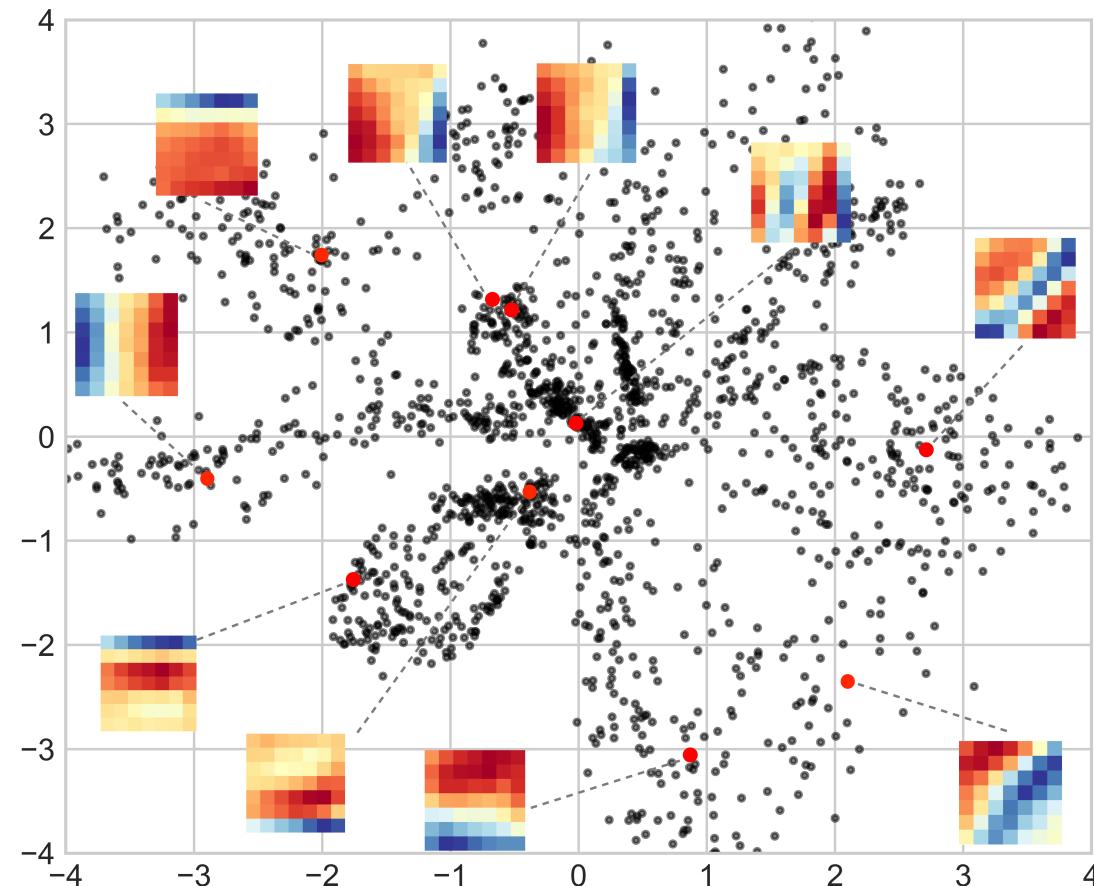
MNIST



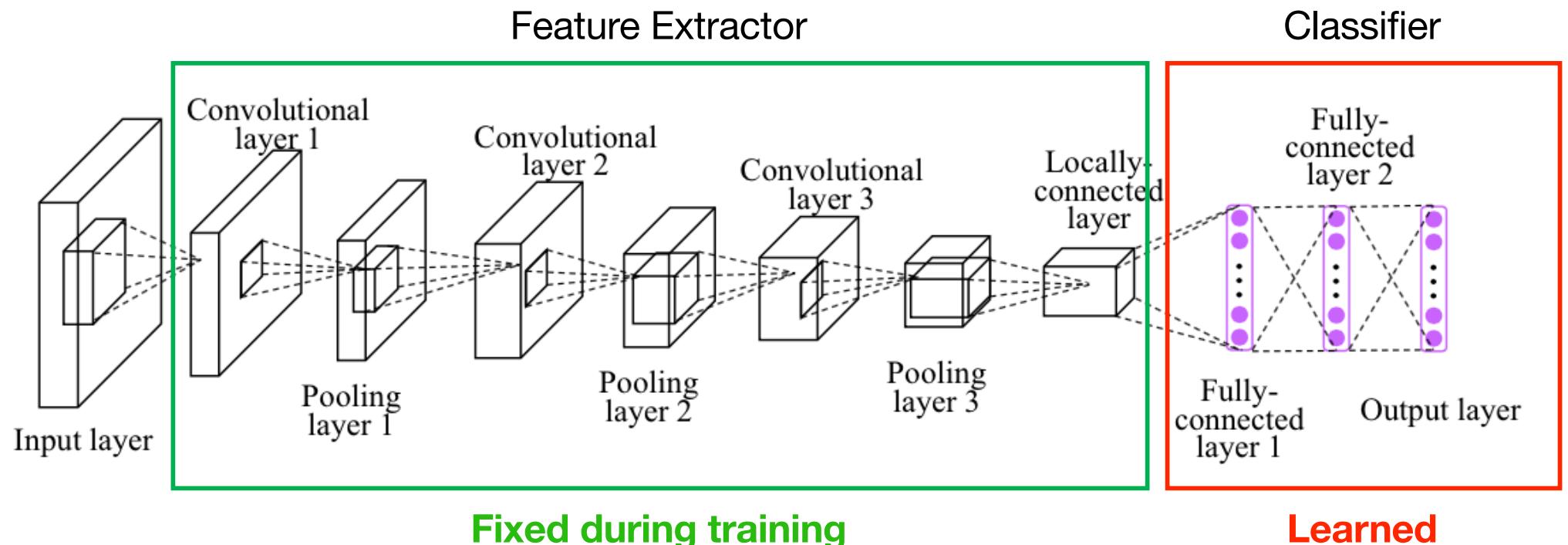
CIFAR-10



DWP Visualization

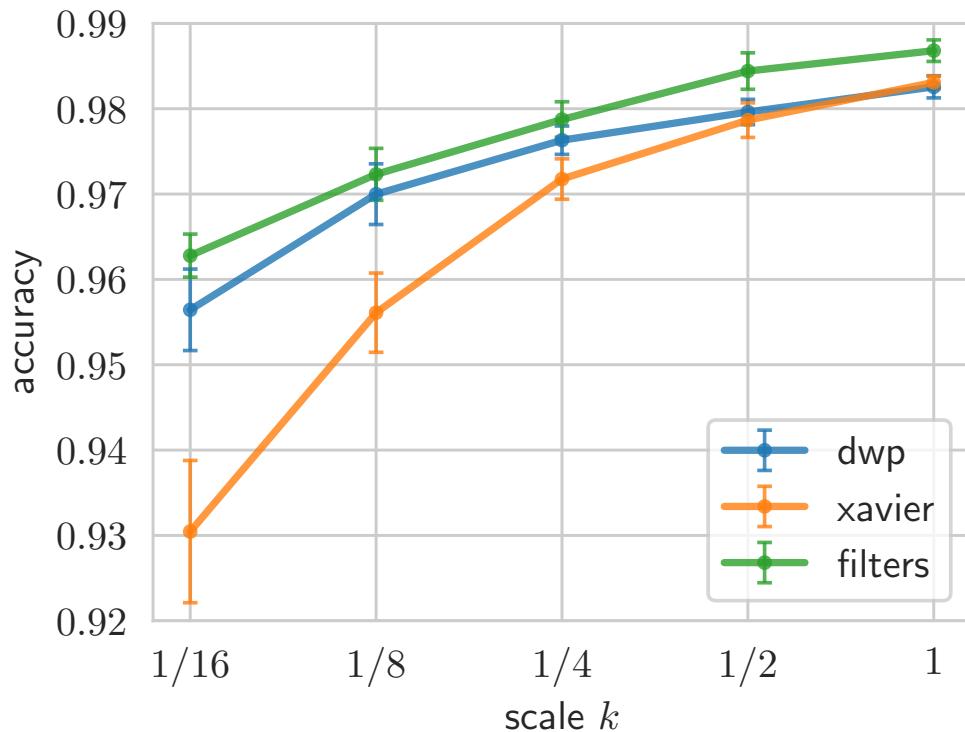


Random Feature Extraction

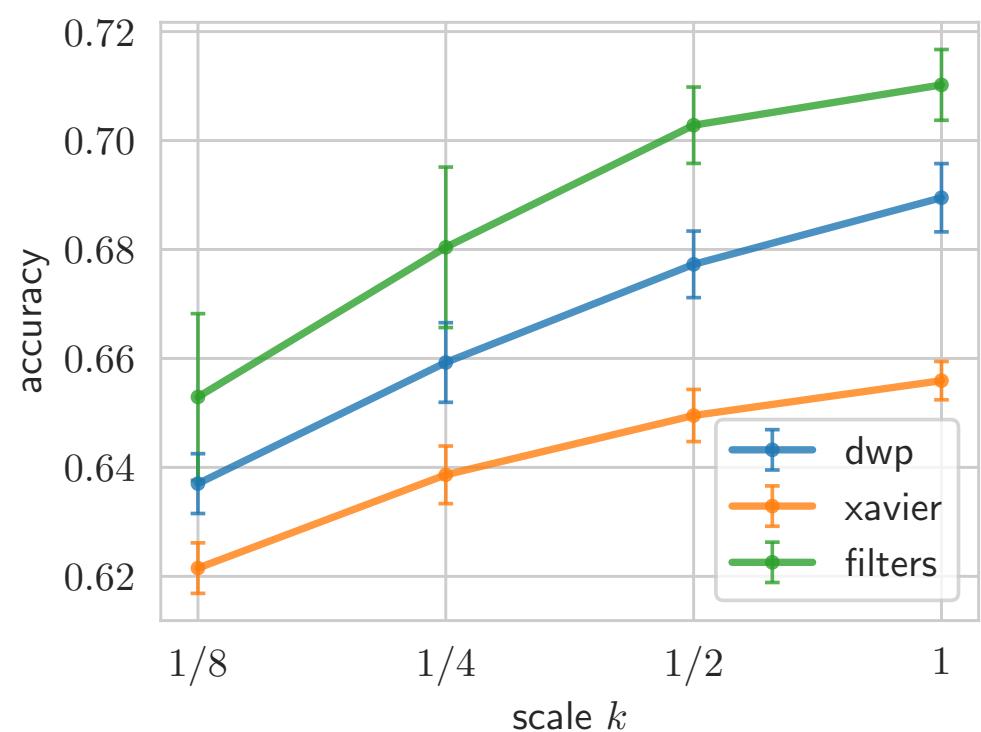


Random Feature Extraction

MNIST

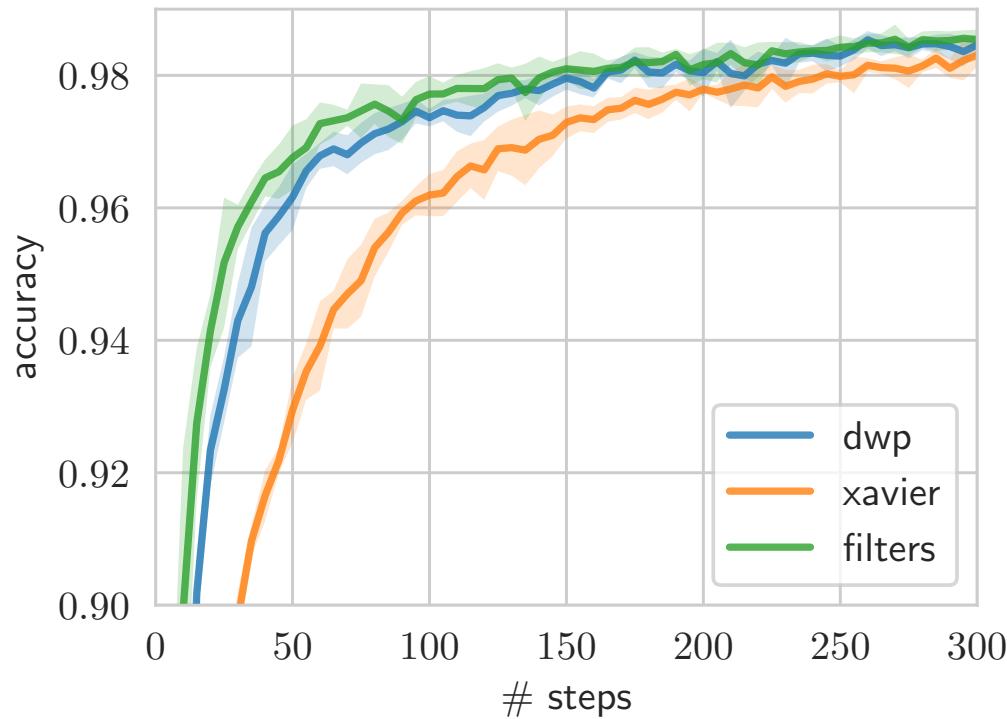


CIFAR-10

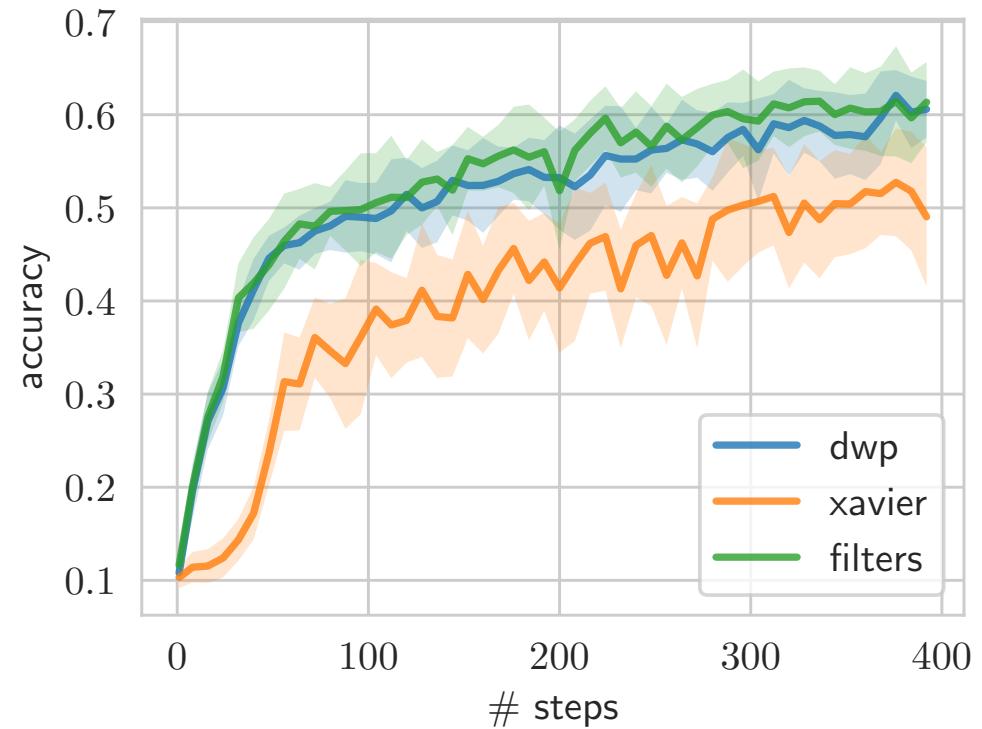


Fast Convergence

MNIST



CIFAR-10



Deep Weight Prior

- Propose the expressive prior distribution capable to encode interdependencies of the learned convolutional filters.
- Develop the method for variational inference with the proposed implicit prior distribution.