

[Paper review 18]

Variational Dropout Sparsifies Deep Neural Networks

(Dmitry Molchannov, et al., 2017)

[Contents]

- 0. Abstract
- 1. Introduction
- 2. Related Research
- 3. Dropout as Bayesian Approximation
- 4. Obtaining Model Uncertainty

0. Abstract

Extend Variational Dropout to the case when dropout rates are unbounded

Propose a way to reduce the variance of the gradient estimator

1. Introduction

Dropout

- Binary Dropout (Hinton et al., 2012)
- Gaussian Dropout (Srivastava et al, 2014)
 - (multiplies the outputs of the neurons by Gaussian random noise)
- Dropout rates are usually optimized by grid-search
 - (To avoid exponential complexity, dropout rates are usually shared for all layers)
- can be seen as a Bayesian regularization (Gal & Ghahramani, 2015)

Instead of injecting noise, Sparsity!

- inducing sparsity during training DNN leads regularization (Han et al., 2015a)
- Sparse Bayesian Learning (Tipping, 2001)
 - (provies framework for training of sparse models)

This paper

- 1) study Variational Dropout (Kingma et al, 2015)

where each weight of a model has its own individual dropout rate

- 2) propose Sparse Variational Dropout
"extends VD" to all possible values of drop out rates ($= \alpha$)
(to do this, provide a new approximation of KL-divergence term in VD objective)
- 3) propose a way to reduce variance of stochastic gradient estimator
→ leads to faster convergence

2. Related Work

너무 많아..생략

논문 참조

3. Preliminaries

3.1 Bayesian Inference

HOW to minimize $D_{KL}(q_\phi(w)||p(w|\mathcal{D}))$?

Maximize ELBO = (1) Expected Log-likelihood - (2) KL-divergence

ELBO : $\mathcal{L}(\phi) = L_{\mathcal{D}}(\phi) - D_{KL}(q_\phi(w)||p(w)) \rightarrow \max_{\phi \in \Phi}$

- (1) Expected Log-likelihood : $L_{\mathcal{D}}(\phi) = \sum_{n=1}^N \mathbb{E}_{q_\phi(w)} [\log p(y_n | x_n, w)]$
- (2) KL-divergence : $D_{KL}(q_\phi(w)||p(w))$

3.2 Stochastic Variational Inference

(a) Reparameterization Trick (Kingma & Welling, 2013)

- obtain unbiased differentiable minibatch-based MC estimator of expected log-likelihood
(that is, find $\nabla_\phi L_{\mathcal{D}}(q_\phi)$)
- trick : decompose into (1) deterministic & (2) stochastic part
 $w = f(\phi, \epsilon)$ where $\epsilon \sim p(\epsilon)$
- number of data in one mini-batch : M

$$\mathcal{L}(\phi) \simeq \mathcal{L}^{SGVB}(\phi) = L_{\mathcal{D}}^{SGVB}(\phi) - D_{KL}(q_\phi(w)||p(w))$$

$$L_{\mathcal{D}}(\phi) \simeq L_{\mathcal{D}}^{SGVB}(\phi) = \frac{N}{M} \sum_{m=1}^M \log p(\tilde{y}_m | \tilde{x}_m, f(\phi, \epsilon_m))$$

$$\nabla_\phi L_{\mathcal{D}}(\phi) \simeq \frac{N}{M} \sum_{m=1}^M \nabla_\phi \log p(\tilde{y}_m | \tilde{x}_m, f(\phi, \epsilon_m))$$

(b) Local Reparameterization Trick (Kingma et al., 2015)

- sample separate weight matrices for each data-point inside mini-batch
- done efficiently by moving the noise from "weights" to "activation"

3.3 Variational Dropout

$B = (A \odot \Xi)W$, with $\xi_{mi} \sim p(\xi)$ putting noise on INPUT

Bernoulli(Binary) Dropout

- Hinton et al., 2012
- $\xi_{mi} \sim \text{Bernoulli}(1 - p)$

Gaussian Dropout with continuous noise

- Srivastava et al, 2014
- $\xi_{mi} \sim \mathcal{N}(1, \alpha = \frac{p}{1-p})$
- continuous noise is better than discrete noise
(multiplying the inputs by Gaussian noise = putting Gaussian noise on the weights)
- can be used to obtain posterior distribution over model's weight! (Wang & Manning, 2013), (Kingma et al., 2015)
($\xi_{ij} \sim \mathcal{N}(1, \alpha) = \text{sampling } w_{ij} \text{ from } q(w_{ij} | \theta_{ij}, \alpha) = \mathcal{N}(w_{ij} | \theta_{ij}, \alpha\theta_{ij}^2) .)$
(Then, $w_{ij} = \theta_{ij}\xi_{ij} = \theta_{ij}(1 + \sqrt{\alpha}\epsilon_{ij}) \sim \mathcal{N}(w_{ij} | \theta_{ij}, \alpha\theta_{ij}^2)$ where $\epsilon_{ij} \sim \mathcal{N}(0, 1)$)

Variational Dropout

- (use reparam trick + draw single sample $W \sim q(W | \theta, \alpha)$)
→ Gaussian dropout = stochastic optimization of exxpected log likelihood
- VD extends this technique!
use $q(W | \theta, \alpha)$ as an approximate posterior with special prior,
 $p(\log|w_{ij}|) = \text{const} \Leftrightarrow p(|w_{ij}|) \propto \frac{1}{|w_{ij}|}$

GD Training = VD Training (when α is fixed)

However, VD provides a way to train dropout rate α by optimizing the ELBO

4. Sparse Variational Dropout

difficulties in training the model with large values of α

→ have considered the case of $\alpha \leq 1$ ($\Leftrightarrow p \leq 0.5$ in binary dropout)

High dropout rate $\alpha_{ij} \rightarrow +\infty = p = 1$

(meaning : corresponding weight is always ignored & can be removed)

4.1 Additive Noise Reparameterization

$$\frac{\partial \mathcal{L}^{SGVB}}{\partial \theta_{ij}} = \frac{\partial \mathcal{L}^{SGVB}}{\partial w_{ij}} \cdot \frac{\partial w_{ij}}{\partial \theta_{ij}} = (1) \times (2)$$

(2) is very noisy if α_{ij} is large.

$$w_{ij} = \theta_{ij} (1 + \sqrt{\alpha_{ij}} \cdot \epsilon_{ij})$$

$$\frac{\partial w_{ij}}{\partial \theta_{ij}} = 1 + \sqrt{\alpha_{ij}} \cdot \epsilon_{ij}, \text{ where } \epsilon_{ij} \sim \mathcal{N}(0, 1)$$

How to reduce variance when α_{ij} is large ?

replace multiplicative noise term $1 + \sqrt{\alpha_{ij}} \cdot \epsilon_{ij}$... with $\sigma_{ij} \cdot \epsilon_{ij}$,

(where $\sigma_{ij}^2 = \alpha_{ij} \theta_{ij}^2$)

$$\begin{aligned} w_{ij} &= \theta_{ij} (1 + \sqrt{\alpha_{ij}} \cdot \epsilon_{ij}) \\ &= \theta_{ij} + \sigma_{ij} \cdot \epsilon_{ij} \end{aligned}$$

Thus, $\frac{\partial w_{ij}}{\partial \theta_{ij}} = 1$, $\epsilon_{ij} \sim \mathcal{N}(0, 1)$

(has no injection noise!)

avoid the problem of large gradient variance!

can train the model within the full range of $\alpha_{ij} \in (0, +\infty)$

4.2. Approximation of the KL Divergence

full KL-divergence term in ELBO

$$D_{KL}(q(W | \theta, \alpha) \| p(W)) = \sum_{ij} D_{KL}(q(w_{ij} | \theta_{ij}, \alpha_{ij}) \| p(w_{ij}))$$

log-scale uniform prior distribution is an improper prior

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

Term above is intractable in VD

need to be sampled & approximated

$$\begin{aligned} -D_{KL}(q(w_{ij} | \theta_{ij}, \alpha_{ij}) \| p(w_{ij})) &\approx \approx k_1 \sigma(k_2 + k_3 \log \alpha_{ij}) - 0.5 \log(1 + \alpha_{ij}^{-1}) + C \\ k_1 &= 0.63576 \quad k_2 = 1.87320 \quad k_3 = 1.48695 \end{aligned}$$

