

## [ Paper review 13 ]

---

# Bayesian Uncertainty Estimation for Batch Normalized Deep Networks

---

( Mattias Teye, et al , 2018 )

## [ Contents ]

---

- 0. Abstract
- 1. Introduction
- 2. Related Works
- 3. Method
  - 1. Bayesian Modeling
  - 2. Batch Normalized Deep Nets as Bayesian Modeling
  - 3. Prior  $p(\omega)$
  - 4. Predictive Uncertainty in Batch Normalized Deep Nets

## 0. Abstract

---

show that "**BN(batch normalization) = approximate inference in Bayesian models**"

- allow us to make estimate of "model uncertainty" using conventional architecture, without modifications to the network!

## 1. Introduction

---

In this work, focus on estimating "**predictive uncertainties in DNN**"

Previous works

1) Drop out (Gal & Ghahramani, 2015)

- any network trained with dropout is an approximate Bayesian Model
- uncertainty estimates can be obtained by "computing the variance of multiple predictions" with different dropout masks
- technique called "MCDO" ( Monte Carlo Dropout )
  - ( can be applied to any pre-trained networks with dropout layers )
  - ( uncertainty estimation for free! )

## 2) Batch normalization

- ability to stabilize learning with improved generalization
- mini-batch statistics depend on randomly selected batch members
  - using this stochasticity, this paper shows that "using BN can be cast as an approximate Bayesian Inference"
  - MCBN ( Monte Carlo Batch Normalization )

# 2. Related Works

---

## Bayesian models for modeling uncertainty

- Gaussian process for infinite parameters ( Neal, 1995 )
- Bayesian NN ( Mackay, 1992 )
- Variational Inference ( Hinton & Van Camp, 1993 ) ( Kingma & Welling, 2014 )
- Probabilistic Backpropagation (PBP) ( Graves, 2011 )
- Factorized posterior via Expectation Propagation ( Hernandez-Lobato & Adams, 2015)
- Deep GP ( Bui et al., 2016 )
- Bayesian Hypernetworks ( Kruger et al., 2017 )
- Multiplicative Normalizing Flows (MNF) ( Louizos & Welling, 2017 )

They all require "MODIFICATION to the architecture"

- Network trained with dropout implicitly performs the VI object
  - Thus, any network trained with dropout can be treated as approximate Bayesian Model ( Gal & Ghahramani, 2015 )
  - ( By making multiple predictions → get mean & variance of them )

# 3. Method

---

## 3.1 Bayesian Modeling

---

deterministic model :  $\hat{y} = \arg \max_y f_{\omega}(\mathbf{x}, \mathbf{y})$

probabilistic model :  $\hat{y} = \arg \max_y f_{\omega}(\mathbf{x}, \mathbf{y}) = \arg \max_y p(\mathbf{y} | \mathbf{x}, \omega)$

- posterior distribution :  $p(\omega | \mathbf{D})$
- probabilistic prediction :  $p(\mathbf{y} | \mathbf{x}, \mathbf{D}) = \int f_{\omega}(\mathbf{x}, \mathbf{y}) p(\omega | \mathbf{D}) d\omega$

## Variational Approximation (VA)

- learn  $q_{\theta}(\omega)$  that minimizes  $\text{KL}(q_{\theta}(\omega) || p(\omega | \mathbf{D}))$
- minimizing  $\text{KL}(q_{\theta}(\omega) || p(\omega | \mathbf{D}))$ 
  - = maximizing ELBO

= minimizing negative ELBO (=

$$\mathcal{L}_{\text{VA}}(\theta) := - \sum_{i=1}^N \int q_{\theta}(\omega) \ln f_{\omega}(\mathbf{x}_i, \mathbf{y}_i) d\omega + \text{KL}(q_{\theta}(\omega) \| p(\omega))$$

= (by MC approximation) minimizing  $\hat{\mathcal{L}}_{\text{VA}}(\theta) := -\frac{N}{M} \sum_{i=1}^M \ln f_{\omega_i}(\mathbf{x}_i, \mathbf{y}_i) + \text{KL}(q_{\theta}(\omega) \| p(\omega))$

(where  $M$  is the size of mini-batch (ex. 64, 128, 256 ...))

$$\hat{\mathcal{L}}_{\text{VA}}(\theta) := -\frac{N}{M} \sum_{i=1}^M \ln f_{\omega_i}(\mathbf{x}_i, \mathbf{y}_i) + \text{KL}(q_{\theta}(\omega) \| p(\omega))$$

- (1) data likelihood :  $-\frac{N}{M} \sum_{i=1}^M \ln f_{\omega_i}(\mathbf{x}_i, \mathbf{y}_i)$
- (2) divergence of the prior w.r.t approximated posterior :  $\text{KL}(q_{\theta}(\omega) \| p(\omega))$

## 3.2 Batch Normalized Deep Nets as Bayesian Modeling

---

inference function :  $f_{\omega}(\mathbf{x}) = W^L a(W^{L-1} \dots a(W^2 a(W^1 \mathbf{x}))$

- $a(\cdot)$  : element-wise non linearity function
- $W^l$  : weight vector at layer  $l$
- $x^l$  : input to layer  $l$
- $h^l = W^l x^l$

### Batch Normalization (BN)

- def) unit-wise operation as below

(standard the distribution of each "unit's input")

$$\hat{h}^u = \frac{h^u - \mathbb{E}[h^u]}{\sqrt{\text{Var}[h^u]}}$$

- during...

1) training : use "mini-batch" (thus, estimated mean & variance on minibatch  $B$  is used)

2) evaluation : use "all training data"

→ therefore, inference at training time for a sample  $x$  is a stochastic process!

(depends on the samples of the mini-batch)

## Loss Function and Optimization

training NN with "mini-batch optimization"

= minimizing  $\mathcal{L}_{\text{RR}}(\omega) := \frac{1}{M} \sum_{i=1}^M l(\hat{\mathbf{y}}_i, \mathbf{y}_i) + \Omega(\omega)$  (regularized risk minimization)

$$\mathcal{L}_{\text{RR}}(\omega) := \frac{1}{M} \sum_{i=1}^M l(\hat{\mathbf{y}}_i, \mathbf{y}_i) + \Omega(\omega)$$

- (1) empirical loss :  $\frac{1}{M} \sum_{i=1}^M l(\hat{\mathbf{y}}_i, \mathbf{y}_i)$
- (2) regularization :  $\Omega(\omega)$

if we set loss function as cross-entropy or SSE , we can also express  $\mathcal{L}_{RR}$  as below

( = minimizing the negative log likelihood )

$$\mathcal{L}_{RR}(\omega) := -\frac{1}{M\tau} \sum_{i=1}^M \ln f_{\omega}(\mathbf{x}_i, \mathbf{y}_i) + \Omega(\omega) \quad (\tau = 1 \text{ for classification})$$

with BN, parameters :  $\{\mathbf{W}^{1:L}, \gamma^{1:L}, \boldsymbol{\beta}^{1:L}, \boldsymbol{\mu}_B^{1:L}, \boldsymbol{\sigma}_B^{1:L}\}$

- learnable params :  $\theta = \{\mathbf{W}^{1:L}, \gamma^{1:L}, \boldsymbol{\beta}^{1:L}\}$
- stochastic params :  $\omega = \{\boldsymbol{\mu}_B^{1:L}, \boldsymbol{\sigma}_B^{1:L}\}$ ,

$$\begin{aligned} \mathcal{L}_{RR}(\omega) &:= -\frac{1}{M\tau} \sum_{i=1}^M \ln f_{\omega}(\mathbf{x}_i, \mathbf{y}_i) + \Omega(\omega) \\ &= -\frac{1}{M\tau} \sum_{i=1}^M \ln f_{\{\theta, \hat{\omega}_i\}}(\mathbf{x}_i, \mathbf{y}_i) + \Omega(\theta) \end{aligned}$$

( where  $\hat{\omega}_i$  : mean & variance for sample  $i$ 's mini-batch )

(  $\hat{\omega}_i$  needs to be i.i.d for training data, but for large number of epochs, it converges to i.i.d cases )

We can estimate uncertainty of predictions by using "**Inherent stochasticity of BN**"

### 3.3 Prior $p(\omega)$

---

for VA & BN to be same....  $\frac{\partial}{\partial \theta}$  of (eq 1) and (eq 2) should be equivalent up to a scaling factor

- (eq 1)  $\hat{\mathcal{L}}_{VA}(\theta) = -\frac{N}{M} \sum_{i=1}^M \ln f_{\omega_i}(\mathbf{x}_i, \mathbf{y}_i) + \text{KL}(q_{\theta}(\omega) \| p(\omega))$
- (eq 2)  $\mathcal{L}_{RR}(\omega) = -\frac{1}{M\tau} \sum_{i=1}^M \ln f_{\{\theta, \hat{\omega}_i\}}(\mathbf{x}_i, \mathbf{y}_i) + \Omega(\theta)$

That is,  $\frac{\partial}{\partial \theta} \text{KL}(q_{\theta}(\omega) \| p(\omega)) = N\tau \frac{\partial}{\partial \theta} \Omega(\theta)$

How to satisfy the condition above?

Solution (1)

- let the prior  $p(\omega)$  imply the regularization term  $\Omega(\theta)$
- in (eq 1), contribution of  $\text{KL}(q_{\theta}(\omega) \| p(\omega))$  to  $\hat{\mathcal{L}}_{VA}$  is "inversly scaled with  $N$ "  
( that is, as  $N \rightarrow \infty$ , NO regularization )

Solution (2)

- let the regularization term  $\Omega(\theta)$  imply the prior  $p(\omega)$
- ex) L-2 regularization :  $\Omega(\theta) = \lambda \sum_{l=1:L} \|W^l\|^2$

### 3.4 Predictive Uncertainty in Batch Normalized Deep Nets

---

approximate predictive distribution :  $p^*(\mathbf{y} \mid \mathbf{x}, \mathbf{D}) := \int f_{\omega}(\mathbf{x}, \mathbf{y}) q_{\theta}(\omega) d\omega$

by Dropout as Bayesian Inference (Gal, 2016)

- mean :  $\mathbb{E}_{p^*}[\mathbf{y}] \approx \frac{1}{T} \sum_{i=1}^T f_{\hat{\omega}_i}(\mathbf{x})$
- covariance :  $\text{Cov}_{p^*}[\mathbf{y}] \approx \tau^{-1} \mathbf{I} + \frac{1}{T} \sum_{i=1}^T f_{\hat{\omega}_i}(\mathbf{x})^{\top} f_{\hat{\omega}_i}(\mathbf{x}) - \mathbb{E}_{p^*}[\mathbf{y}]^{\top} \mathbb{E}_{p^*}[\mathbf{y}]$   
( where  $\hat{\omega}_i$  corresponds to sampling the net's stochastic params  $\omega = \{\mu_B^{1:L}, \sigma_B^{1:L}\}$  )  
( Sampling  $\hat{\omega}_i$  involves sampling a batch  $B$  from training set )

## Algorithm summary

network is trained just as a regular BN network!

But, instead of replacing  $\omega = \{\mu_B^{1:L}, \sigma_B^{1:L}\}$  with population values from  $D$ ,

we update these parameters stochastically, once for each forward pass

---

### Algorithm 1 MCBN Algorithm

---

**Input:** sample  $x$ , number of inferences  $T$ , batchsize  $b$

**Output:** mean prediction  $\hat{y}$ , predictive uncertainty  $\sigma^2$

- 1:  $\mathbf{y} = \{\}$
  - 2: **loop** for  $T$  iterations
  - 3:  $B \sim D$  // mini batch
  - 4:  $\hat{\omega} = \{\mu_B, \sigma_B\}$  // mini batch mean and variance
  - 5:  $\mathbf{y} = \mathbf{y} \cup f_{\hat{\omega}}(x)$
  - 6: **end loop**
  - 7:  $\hat{y} = \mathbb{E}[\mathbf{y}]$
  - 8:  $\sigma^2 = \text{Cov}[\mathbf{y}] + \tau^{-1} \mathbf{I}$  // for regression
-