

1 Notations

1.1 Model & Data

Let $f(x; \theta)$ a function of x with parameter θ . Let $y = f(x; \theta)$ an observable, thus the observed value obeys a Gaussian distribution. Let D denotes a set of observations, $D := \{(x_i, y_i, \sigma_i) : i = 1, \dots, N_D\}$, wherein x_i is the i th input, y_i its observed value, and σ_i the observational error of y_i . We may employ mini-batch technique, thus denote $D_m := \{(x_i, y_i, \sigma_i) : i = 1, \dots, N_m\} \subset D$ as a mini-batch, with batch-size $N_m \leq N_D$.

2 Bayesian

2.1 Prior-Posterior Iteration

2.2 Bayesian as Information Encoder

Comparing with the traditional method, what is the advantage of Bayesian way? The answer is, it encodes more information of data into model. Indeed, it does not encodes the value of peak of the posterior only, as traditional method does, but also much more information on the posterior. XXX

3 Neural Network for Posterior (nn4post)

3.1 The Model

Suppose we have some prior on θ , $p(\theta)$, we gain the unnormalized posterior $p(D|\theta) p(\theta)$. With D arbitrarily given, this unnormalized posterior is a function of θ , denoted by $p(\theta; D)$ ^{1,2}.

We we are going to do is fit this $p(\theta; D)$ by ANN for any given D . To do so, we have to assume that $\text{supp}\{p(\theta; D)\} = \mathbb{R}^d$ for some $d \in \mathbb{N}^+$ (i.e. has no compact support) but decrease exponentially fast as $\|\theta\| \rightarrow +\infty$. With this assumption, $\ln p(\theta; D)$ is well-defined. For ANN, we propose using Gaussian function as the activation-function. Thus, we have the fitting function

$$q(\theta; a, \mu, \zeta) = \sum_{i=1}^{N_c} c_i(a) \left\{ \prod_{\alpha=1}^d \Phi(\theta_\alpha - \mu_{i\alpha}, \sigma(\zeta_{i\alpha})) \right\},$$

where

$$\begin{aligned} c_i(a) &= \frac{\exp(a_i)}{\sum_{j=1}^N \exp(a_j)} = \text{softmax}(i; a); \\ \sigma(\zeta_{i\alpha}) &= \ln(1 + \exp(\zeta_{i\alpha})), \end{aligned}$$

and $a_i, \mu_{i\alpha}, \zeta_{i\alpha} \in \mathbb{R}$ for $\forall i, \forall \alpha$ and

$$\Phi(x; \mu, \sigma) := \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

being the Gaussian PDF. The introduction of ζ is for numerical consideration, see below.

1. This is why we use “;” instead of “, ”, indicating that D has been (arbitrarily) given and fixed.

2. The normalized posterior $p(\theta|D) = p(D|\theta) p(\theta)/p(D) = p(\theta; D)/p(D)$, by Bayes's rule.

3.1.1 Numerical Consideration

If, in q , we regard w , μ , and σ as independent variables, then the only singularity appears at $\sigma=0$. Indeed, σ appears in Φ (as well as the derivatives of Φ) as denominator only, while others as numerators. However, once doing numerical iterations with a finite step-length of σ , the probability of reaching or even crossing 0 point cannot be surely absent. This is how we may encounter this singularity in practice.

Introducing the ζ is our trick of avoiding this singularity. Precisely, using a singular map that pushes the singularity to infinity solves the singularity. In this case, using $\text{softplus}(\cdot)$ that pushes $\sigma=0$ to $\zeta \rightarrow -\infty$, so that, with finite steps of iteration, singularity (at $-\infty$) cannot be reached.

This trick (i.e. pushing a singularity to infinity) is the same as in avoiding the horizon-singularity of Schwarzschild solution of black hole.

3.2 Interpretation

3.2.1 As a Mixture Distribution

$q(\theta; a, \mu, \zeta)$ has a probabilistic interpretation. $\prod_{j=1}^d \Phi(\theta_j - \mu_{ij}, \sigma(\zeta_{ij}))$ corresponds to multi-dimensional Gaussian distribution (denote \mathcal{N}), with all dimensions independent with each other. The $\{c_i(a)\}$ is a categorical distribution, randomly choosing the Gaussian distributions. Thus $q(\theta; a, \mu, \zeta)$ is a composition: categorical \rightarrow Gaussian. This is the *mixture distribution*.

3.2.2 As a Generalization

This model can also be interpreted as a direct generalization of mean-field variational inference. Indeed, let $N_c = 1$, this model reduces to mean-field variational inference. Remark that mean-field variational inference is a mature algorithm and has been successfully established on many practical applications.

3.2.3 As a Neural Network

3.3 Marginalization

This model can be marginalized easily. This then benefits the transferring of the model components. Precisely, for any dimension-index β given, we can marginalize all other dimensions directly, leaving

$$\begin{aligned} q(\theta_\beta; a, \mu, \zeta) &= \prod_{\forall \gamma \neq \beta} \int d\theta_\gamma \sum_{i=1}^{N_c} c_i(a) \left\{ \prod_{\alpha=1}^d \Phi(\theta_\alpha; \mu_{i\alpha}, \sigma(\zeta_{i\alpha})) \right\} \\ &= \sum_{i=1}^{N_c} c_i(a) \Phi(\theta_\beta; \mu_{i\beta}, \sigma(\zeta_{i\beta})), \end{aligned}$$

where employed the normalization of Φ .

3.4 Loss-Function

We use “evidence of lower bound” (ELBO) as loss. It is ensured to have a unique global minimal, at which $p(\theta; D) = q(\theta; a, \mu, \zeta)$.

$$\begin{aligned} \text{ELBO}(a, \mu, \zeta) &:= \mathbb{E}_{\theta \sim q(\theta; a, \mu, \zeta)} [\ln p(\theta; D) - \ln q(\theta; a, \mu, \zeta)] \\ &\approx \left(\frac{1}{n} \sum_{\theta(s)} \right) \{ \ln p(\theta(s); D) - \ln q(\theta(s); a, \mu, \zeta) \}, \end{aligned}$$

where $\{\theta_{(s)}; s = 1, \dots, n\}$ is sampled from $q(\theta; a, \mu, \zeta)$ as a distribution. Since there's no compact support for both $p(\theta; D)$ and $q(\theta; a, \mu, \zeta)$, ELBO is well-defined, as the loss-function (or say loss-function, performance, etc) of the fitting.

4 Stochastic Optimization

4.1 Difference between Bayesian and Traditional Methods

Suppose, instead of use the whole dataset, we employ mini-batch technique. Since all data are independent, if suppose that D_m is unbiased in D , then we have,

$$\ln p(D|\theta) = \sum_D p((x_i, y_i, \sigma_i)|\theta) \approx \frac{N_D}{N_m} \sum_{D_m} p((x_i, y_i, \sigma_i)|\theta) = \frac{N_D}{N_m} \ln p(D_m|\theta).$$

Then,

$$\ln p(\theta; D) = \ln p(D|\theta) + \ln p(\theta) = \frac{N_D}{N_m} \ln p(D_m|\theta) + \ln p(\theta),$$

thus as previous

$$\ln p(\theta; D) = \frac{N_D}{N_m} \sum_{(x_i, y_i, \sigma_i) \in D_m} \left\{ -\frac{1}{2} \ln(2\pi\sigma_i^2) - \frac{1}{2} \left(\frac{y_i - f(x_i; \theta)}{\sigma_i} \right)^2 \right\} + \ln p(\theta).$$

In this we meet one of the main differences between the Bayesian and the traditional. In the traditional method, N_D does not matters in training, being absent in the optimizer. However, in Bayesian, the number of data that are employed is encoded into Bayesian model, and has to, since the greater number of data gives more confidence. So, while using stochastic optimization in Bayesian mode, the factor N_D/N_m of likelihood has to be taken into account. We have to know how many data we actually have, thus how confident we are.

5 ADVI

Automatic differentiation variational inference (ADVI)³ has the advantage that the variance of its Monte Carlo integral is orderly smaller than that of black box variational inference (i.e. optimization directly using ELBO without further reparameterization).

Precisely, let \mathbb{E} for mean value, \mathbb{H} for shannon entropy, Φ for Gaussian, and $\sigma(\cdot)$ for softplus function. By

$$q(\theta; a, \mu, \zeta) = \sum_{i=1}^{N_c} w_i(a) \Phi(\theta; \mu_i, \sigma(\zeta_i)),$$

we have

$$\begin{aligned} \text{ELBO} &= \mathbb{E}_{q(\theta; a, \mu, \zeta)} [\ln p(\theta; D)] + \mathbb{H}[q(\theta; a, \mu, \zeta)] \\ &= \sum_i^{N_c} w_i(a) \mathbb{E}_{\Phi_i(\theta; \mu_i, \sigma(\zeta_i))} [\ln p(\theta; D)] + \mathbb{H}[q(\theta; a, \mu, \zeta)] \\ &=: E_1 + E_2. \end{aligned}$$

(E_2 is analytic and independent of $p(\theta; D)$, so we leave it for later.) Then, for $\forall i = 1, \dots, N_c$, $\forall \alpha = 1, \dots, N_d$, let

$$\eta_\alpha := \frac{\theta_\alpha - \mu_{i\alpha}}{\sigma(\zeta_{i\alpha})},$$

³. See, Kucukelbir, et al, 2016.

we have $\theta_\alpha = \sigma(\zeta_{i\alpha}) \eta_\alpha + \mu_{i\alpha}$ ($\theta = \sigma(\zeta_i) \eta + \mu_i$ if hides the α index). So, for any i -components in the E_1 of ELBO, we transform

$$\mathbb{E}_{\Phi(\theta; \mu_i, \sigma(\zeta_i))}[\ln p(\theta; D)] = \mathbb{E}_{\Phi(\eta; 0, 1)}[\ln p(\sigma(\zeta_i) \eta + \mu_i; D)].$$

Thus, we have derivatives

$$\begin{aligned}\frac{\partial E_1}{\partial \mu_{i\alpha}} &= w_i(a) \mathbb{E}_{\Phi(\eta; 0, 1)}[\nabla_\alpha \ln p(\sigma(\zeta_i) \eta + \mu_i; D)]; \\ \frac{\partial E_1}{\partial \zeta_{i\alpha}} &= w_i(a) \mathbb{E}_{\Phi(\eta; 0, 1)}\left[\nabla_\alpha \ln p(\sigma(\zeta_i) \eta + \mu_i; D) \eta_\alpha \frac{\partial \sigma}{\partial \zeta_{i\alpha}}(\zeta_i)\right],\end{aligned}$$

where $\nabla_\alpha \ln p := \partial \ln p / \partial \theta_\alpha$ as the formal derivative. So, for these two, value of $\ln p(\theta; D)$ is regardless.

6 Deep Learning

It cannot solve the vanishing gradient problem of deep neural network, since this problem is intrinsic to the posterior of deep neural network. Indeed, the posterior has the shape like $\exp(-x^2/\sigma^2)$ with $\sigma \rightarrow 0$, where x is the variable (argument) of the posterior. It has a sharp peak, located at a tiny area, with all other region extremely flat. The problem of find this peak, or equivalently, finding its tiny area, is intrinsically intactable.

So, even for Bayesian neural network, a layer by layer abstraction along depth cannot be absent.

7 Transfer Learning

8 Why not MCMC?