

Stochastic Natural Gradient Algorithm and Online learning for Latent Dirichlet Allocation

Nanbo Sun

Electrical & Computer Engineering, National University of Singapore, Singapore.
nanbosun@u.nus.edu

Abstract. This report derives in detail a stochastic optimization algorithm — stochastic natural gradient algorithm [1], and apply it to Latent Dirichlet Allocation. We first introduce the concept of stochastic gradient descent and natural gradient. Next, we combine them to create the stochastic natural gradient algorithm by using the information geometry. Finally, we apply the algorithm to Latent Dirichlet Allocation.

1 Stochastic Natural gradient algorithm

Let us introduce these concept by using an example of multilayer perceptron. Each **sample** z is a pair composed of an **input variable** x and a **desired output** y . The output of the network is a function $f(x, w)$ of input variable x and the **weight vector** w . And, there are L samples $\{z_i\}, i = 1 \dots L$.

1.1 Stochastic gradient descent

The performance of the network is measured by a **loss function** $Q(z, w)$. Usually, the loss function is the squared mean error which is defined as following:

$$Q_{mse}(z, w) \triangleq \frac{1}{2}(y - f(x, w))^2 \quad (1)$$

Then, the expected risks and empirical risk can be written as:

$$C(w) \triangleq \mathbf{E}_z Q(z, w) \triangleq \int Q(z, w)p(z)dz \approx \frac{1}{L} \sum_i Q(z_i, w) \quad (2)$$

where $p(z)$ is the *ground truth* distribution of sample z . In the case of the empirical risk, the expectation then reduces to a simple average $1/L \sum_i Q(z_i, w)$.

If we use *gradient descent* algorithm to optimize it, the update equation is:

$$w_{t+1} = w_t - \gamma_t \nabla_w C(w) \approx w_t - \gamma_t \sum_i \nabla_w Q(z_i, w_t) \quad (3)$$

If we use *stochastic gradient descent* algorithm to optimize it, the update equation is:

$$w_{t+1} = w_t - \gamma_t \nabla_w Q(z_t, w_t) \quad (4)$$

The algorithm presents the concept of *stochastic approximation*. It's hoped that Eq.4 behaves like Eq.3.

Only if it satisfies the following 4 conditions [2], we can prove that the stochastic gradient descent is convergent.

$$\forall \epsilon > 0, \quad \inf_{(w-w^*)^2 > \epsilon} (w-w^*)\nabla_w C(w) > 0 \quad (5)$$

The hypothesis simply states that the opposite of the gradient always points toward the minimum w^* .

In order to obtain convergence, the parameter updates $\gamma_t \nabla_w Q(z, w)$ must become smaller and smaller when the parameter vector approaches the optimum w^* . The size of the gradients can be described with a first order expansion of $\nabla_w Q$ in w^* .

$$\mathbf{E}_z(\nabla_w Q(z, w)^2) \approx \mathbf{E}_z(\nabla_w Q(z, w^*)^2) + \mathbf{E}((w-w^*) \dots) < A + B(w-w^*)^2 \quad (6)$$

where $A, B \geq 0$ are constants. When $w \rightarrow w^*$, $\mathbf{E}_z(\nabla_w Q(z, w)^2) \rightarrow A$. Therefore, we must use *decreasing learning rates*

$$\sum \gamma_t^2 < \infty \quad (7)$$

Suppose the learning rates decrease so fast that $\sum \gamma_t = R < \infty$. This would effectively maintain the parameters within a certain radius of their initial value. Therefore

$$\sum \gamma_t = \infty \quad (8)$$

1.2 Natural Gradient

Assume $S = \{\mathbf{w} \in R^n\}$ is a parameter space and $C(w)$ is a function defined in the space. Then, the distance between \mathbf{w} and $\mathbf{w} + d\mathbf{w}$ is defined as following:

$$\begin{aligned} \text{Euclidean space: } \quad \|\mathbf{d}\mathbf{w}\|^2 &= \sum_{i=1}^n (dw_i)^2 \\ \text{Riemannian space: } \quad \|\mathbf{d}\mathbf{w}\|^2 &= \sum_{i,j} g_{ij}(\mathbf{w}) dw_i dw_j \end{aligned}$$

where $G = (g_{ij})$ is the Riemannian metric tensor.

Theorem 1. *The steepest descent direction of $C(w)$ in Riemannian space is [3]*

$$-\tilde{\nabla}C(w) = -G^{-1}(w)\nabla C(w) \quad (9)$$

Proof. Assume $d\mathbf{w} = \epsilon \mathbf{a}$, we want to minimize $C(\mathbf{w} + d\mathbf{w})$.

$$\begin{aligned} \min_{\mathbf{a}} C(\mathbf{w} + d\mathbf{w}) &= C(\mathbf{w}) + \epsilon \nabla C(\mathbf{w}) \mathbf{a} \\ \text{s.t. } \|\mathbf{a}\|^2 &= \sum g_{ij} a_i a_j = 1 \end{aligned}$$

We can use the Lagrangian method to solve this problem

$$\begin{aligned}\frac{\partial}{\partial a_i} \{ \nabla C(\mathbf{w})^T \mathbf{a} - \lambda \mathbf{a}^T G \mathbf{a} \} &= 0 \\ \mathbf{a} &= \frac{1}{2\lambda} G^{-1} \nabla C(\mathbf{w}) \\ \Rightarrow -\tilde{\nabla} C(\mathbf{w}) &= -G^{-1}(\mathbf{w}) \nabla C(\mathbf{w})\end{aligned}$$

1.3 Information Geometry

Assume $y = f(x, w) + \epsilon$, where ϵ is the Gaussian noise.

$$p(z, w) = C \exp\left[-\frac{(y - f(x, w))^2}{2\sigma^2}\right] \quad (10)$$

Because we hope that $p(z, w) \approx p(z)$ which is the *ground truth* distribution of z , we can minimize the following equation

$$\min_w D_{kl}(p(z)||p(z, w)) = \mathbf{E}_z \log \frac{p(z)}{p(z, w)} = \frac{1}{\sigma^2} \mathbf{E}_z Q_{mse}(z, w) + \text{Constant} \quad (11)$$

According to the above equation, we can find that it is the same as minimizing the *mean square error* loss function. We can measure the distance between \mathbf{w} and $\mathbf{w} + d\mathbf{w}$ by following equation

$$D_{kl}(\mathbf{w}||\mathbf{w} + d\mathbf{w}) = D(\mathbf{w}||\mathbf{w} + d\mathbf{w})|_{\mathbf{w}=\mathbf{w}+d\mathbf{w}} \quad (12)$$

$$+ d\mathbf{w} \cdot \frac{d}{d\mathbf{w}} D(\mathbf{w}||\mathbf{w} + d\mathbf{w})|_{\mathbf{w}=\mathbf{w}+d\mathbf{w}} \quad (13)$$

$$+ \frac{1}{2} d\mathbf{w}^T \left[\frac{d^2}{d\mathbf{w}_i d\mathbf{w}_j} D(\mathbf{w}||\mathbf{w} + d\mathbf{w}) \right] |_{\mathbf{w}=\mathbf{w}+d\mathbf{w}} d\mathbf{w} \quad (14)$$

$$= \frac{1}{2} d\mathbf{w}^T \underbrace{\left[\frac{d^2}{d\mathbf{w}_i d\mathbf{w}_j} D(\mathbf{w}||\mathbf{w} + d\mathbf{w}) \right] |_{\mathbf{w}=\mathbf{w}+d\mathbf{w}}}_{\text{Fisher Information}} d\mathbf{w} \quad (15)$$

It is obvious that the Fisher Information can be regarded as a Riemannian metric tensor G .

1.4 Stochastic Natural Gradient Algorithm

Instead of using the gradient, we can use the natural gradient to substitute it. Therefore, we can get the following equation

$$w_{t+1} = w_t - \gamma_t G^{-1}(w_t) \nabla_w Q(z_t, w_t) \quad (16)$$

where $G = \int -(\nabla_w^2 \log p(z, w)) p(z, w) dw$.

2 Online Learning for LDA [4]

Now, let's apply the stochastic gradient descent algorithm to Latent Dirichlet Allocation. From Eq.(11) in [4], we know that

$$\mathcal{L}(g, \boldsymbol{\lambda}) \triangleq D\mathbf{E}_g[l(n, \gamma(n, \boldsymbol{\lambda}), \phi(n, \boldsymbol{\lambda}), \boldsymbol{\lambda})|\boldsymbol{\lambda}]. \quad (17)$$

Using the stochastic natural gradient descent, we get

$$\boldsymbol{\lambda} \leftarrow \boldsymbol{\lambda} + \rho_t D G^{-1} \nabla_{\boldsymbol{\lambda}} l(n_t, \gamma_t, \phi_t, \boldsymbol{\lambda}) \quad (18)$$

According to Eq.(14), we know that $G^{-1} \nabla_{\boldsymbol{\lambda}} l(n_t, \gamma_t, \phi_t, \boldsymbol{\lambda}) = -\lambda_{kw}/D + \eta/D + n_{tm} \phi_{twk}$. Then, the update equation becomes

$$\boldsymbol{\lambda} \leftarrow \boldsymbol{\lambda} + \rho_t D (-\lambda_{kw}/D + \eta/D + n_{tm} \phi_{twk}) \quad (19)$$

$$\leftarrow \boldsymbol{\lambda} + \rho_t (-\lambda_{kw} + \underbrace{\eta + D n_{tm} \phi_{twk}}_{\tilde{\lambda}_{kw}}) \quad (20)$$

$$\leftarrow (1 - \rho_t) \boldsymbol{\lambda} + \rho_t \tilde{\boldsymbol{\lambda}} \quad (21)$$

The above equation is the same as that in **Algorithm 2**.

References

1. Bottou, Lon, and Noboru Murata. "Stochastic approximations and efficient learning." The Handbook of Brain Theory and Neural Networks, Second edition,. The MIT Press, Cambridge, MA (2002).
2. Bottou, Lon. "Online learning and stochastic approximations." On-line learning in neural networks 17.9 (1998): 25.
3. Amari, Shun-Ichi. "Natural gradient works efficiently in learning." Neural computation 10.2 (1998): 251-276.
4. Hoffman, Matthew, Francis R. Bach, and David M. Blei. "Online learning for latent dirichlet allocation." advances in neural information processing systems. 2010.