

# Progress report

Lukas Zorich

April 28, 2017

## 1 Background

Let  $X = \{x_1, \dots, x_N\}$  a collection of  $N$  points. Now, lets suppose that the features of  $x_j$  are updated, and lets call  $x_j^*$  the updated point. Lastly, lets define  $X_{-j} = \{x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_N\}$ .

The posterior before  $x_j$  moves, and using  $P(\Theta)$  as the prior is

$$P(\Theta | X) \propto P(X | \Theta)P(\Theta)$$

,

where  $P(X | \Theta) = \prod_i^N P(x_i | \Theta)$ .

Now, after  $x_j$  moves to  $x_j^*$ , the posterior is

$$P(\Theta | X_{-j}, x_j^*) \propto P(X_{-j}, x_j^* | \Theta)P(\Theta) \quad (1)$$

$$\propto P(X_{-j} | \Theta)P(x_j^* | \Theta)P(\Theta) \quad (2)$$

$$\propto \frac{P(X | \Theta)}{P(x_j | \Theta)}P(x_j^* | \Theta)P(\Theta) \quad (3)$$

$$\propto \frac{P(x_j^* | \Theta)}{P(x_j | \Theta)}P(X | \Theta)P(\Theta) \quad (4)$$

$$(5)$$

which can be written as,

$$P(\Theta | X_{-j}, x_j^*) \propto \frac{P(x_j^* | \Theta)}{P(x_j | \Theta)}P(\Theta | X) \quad (6)$$

Now, if we instead of considering one data point  $x_j$ , we consider a batch of  $S$  data points  $X_J = \{x_j, x_{j+1}, \dots, x_{j+S-1}, x_{j+S}\}$  that moved from  $x_{j+k}$  to  $x_{j+k}^*$ , for  $k = 0, \dots, S$ . And lets define  $X_{-J} = \{x_1, \dots, x_{j-1}, x_{j+S+1}, \dots, x_N\}$ . Then, Eq. 6 for the batch of  $S$  data point and because the prior  $P(\Theta)$  gets canceled, can be written as,

$$P(\Theta | X_{-J}, X_J^*) \propto \frac{P(\Theta | X_J^*)}{P(\Theta | X_{-J})}P(\Theta | X) \quad (7)$$

Inspired by the work done in [1], we assume that we approximate the posterior using **variational inference**. Also, we assume that  $P(\Theta)$  is an exponential family distribution for  $\Theta$  with sufficient statistic  $T(\Theta)$  and natural parameter  $\xi_0$ . We suppose further that if  $q(\Theta)$  is the approximate posterior obtained using variational inference, then  $q(\Theta)$  is also in the same exponential family with a parameter  $\xi$  such that

$$q(\Theta) \propto \exp(\xi \cdot T(\Theta)). \quad (8)$$

Similar to [1], when we make this assumptions the update in Eq. 7 becomes

$$P(\Theta \mid X_{-J}, X_J^*) \approx \exp([\xi - \xi_J + \xi_J^*] \cdot T(\Theta)) \quad (9)$$

where  $\xi$  is the natural parameter of  $q(\Theta) \approx P(\Theta \mid X)$ , and  $\xi_J$  and  $\xi_J^*$  corresponds to the natural parameter of  $q_J(\Theta) \approx p(\Theta \mid X_J)$  and  $q_J^*(\Theta) \approx p(\Theta \mid X_J^*)$  respectively.

**Using this approach, we can update the posterior when data "moves" without the need to go through the whole dataset, instead we just need to go through the data points that moved to obtain  $\xi_J$  and  $\xi_J^*$ .**

## 2 Application to a simplified version of my model

For trying the proposed approach, I'm starting only with a single GMM (I prefer to start small).

For a GMM, the update would be

$$P(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda} \mid X_{-J}, X_J^*) \approx \frac{q_{-J}^*(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})}{q_{-J}(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \quad (10)$$

$$(11)$$

where

$$q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = q(\boldsymbol{\pi})p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \quad (12)$$

$$= \text{Dir}(\boldsymbol{\pi} \mid \boldsymbol{\alpha})\mathcal{N}(\boldsymbol{\mu} \mid \mathbf{m}, (\beta\boldsymbol{\Lambda})^{-1})\mathcal{W}(\boldsymbol{\Lambda} \mid \mathbf{W}, \nu). \quad (13)$$

### 2.1 Updates

#### 2.1.1 Dirichlet

The natural parameter for the dirichlet is:

$$\xi = \boldsymbol{\alpha} - 1 \quad (14)$$

hence, the update is:

$$\boldsymbol{\alpha}' = \boldsymbol{\alpha} - \boldsymbol{\alpha}_J + \boldsymbol{\alpha}_J^* \quad (15)$$

### 2.1.2 Normal-Wishart

The natural parameter for the Normal-Wishart distribution is:

$$\xi = \begin{bmatrix} \beta \boldsymbol{\mu} \\ \beta \\ \boldsymbol{\Lambda}^{-1} + \beta \boldsymbol{\mu} \boldsymbol{\mu}^T \\ \nu + 2 + p \end{bmatrix} \quad (16)$$

hence, the updates are:

$$\boldsymbol{\mu}' = \frac{1}{\beta'} (\beta \boldsymbol{\mu} - \beta_J \boldsymbol{\mu}_J + \beta_J^* \boldsymbol{\mu}_J^*) \quad (17)$$

$$\beta' = \beta - \beta_J + \beta_J^* \quad (18)$$

$$\boldsymbol{\Lambda}^{-1'} = (\boldsymbol{\Lambda}^{-1} + \beta \boldsymbol{\mu} \boldsymbol{\mu}^T) - (\boldsymbol{\Lambda}_J^{-1} + \beta_J \boldsymbol{\mu}_J \boldsymbol{\mu}_J^T) + (\boldsymbol{\Lambda}_J^{-1*} + \beta_J^* \boldsymbol{\mu}_J^* \boldsymbol{\mu}_J^{T*}) - \beta' \boldsymbol{\mu}' \boldsymbol{\mu}'^T \quad (19)$$

$$\nu' = \nu - \nu_J + \nu_J^* \quad (20)$$

$$(21)$$

## 3 Current progress

During the three weeks since I sent my goals for this months I distributed my time like this:

1. Week 1: Try to implement the online update in PyMC. I tried to implement this, but I kept getting errors. I asked in the PyMC forum for help, but nobody could help me. I think that PyMC is not flexible to do this type of work.
2. Week 2: Decided to move to variational inference, because there are several papers that use variational inference for streaming bayesian models. I started implementing the variational inference GMM.
3. Week 3: Finished implementing the variational inference GMM. Started working in the online update, but I'm having some bugs, so this part is not finished yet.

The current progress is the following:

- Implementing variational bayes GMM: 100%.
- Implementing online update of the posterior: 70%. I have the code written, but I'm having a numerical underflow bug when doing the updates.

## 4 Next steps

Finish implementing this (hopefully this weekend), and depending on the results decide the next steps.

## References

- [1] BRODERICK, T., BOYD, N., WIBISONO, A., WILSON, A. C., AND JORDAN, M. I. Streaming variational bayes. In *NIPS* (2013), C. J. C. Burges, L. Bottou, Z. Ghahramani, and K. Q. Weinberger, Eds., pp. 1727–1735.