

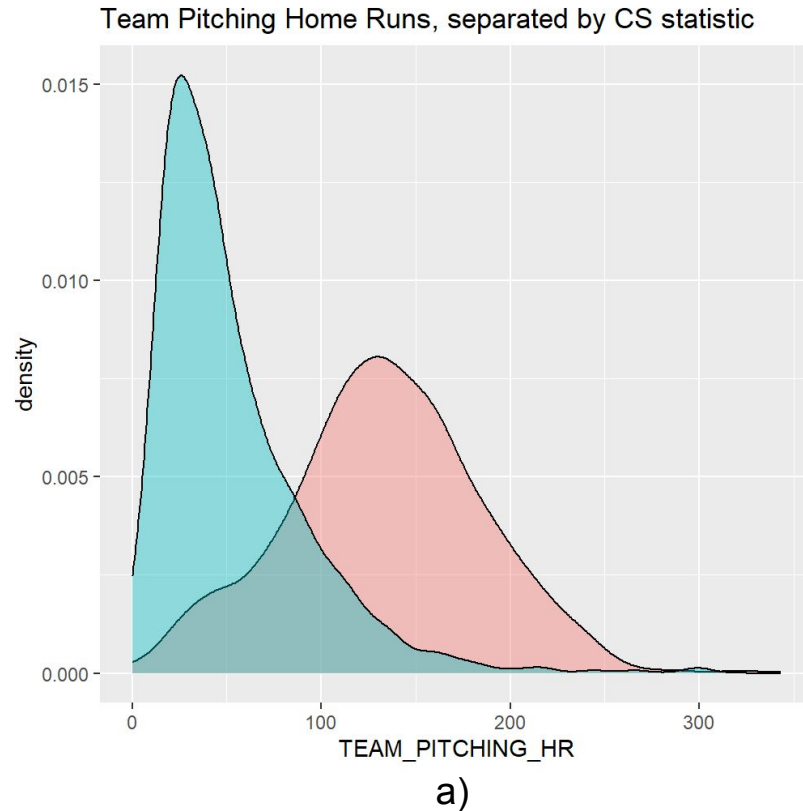
Automated Variational Inference

Daniela Toader & Radu Gaghi

Introduction

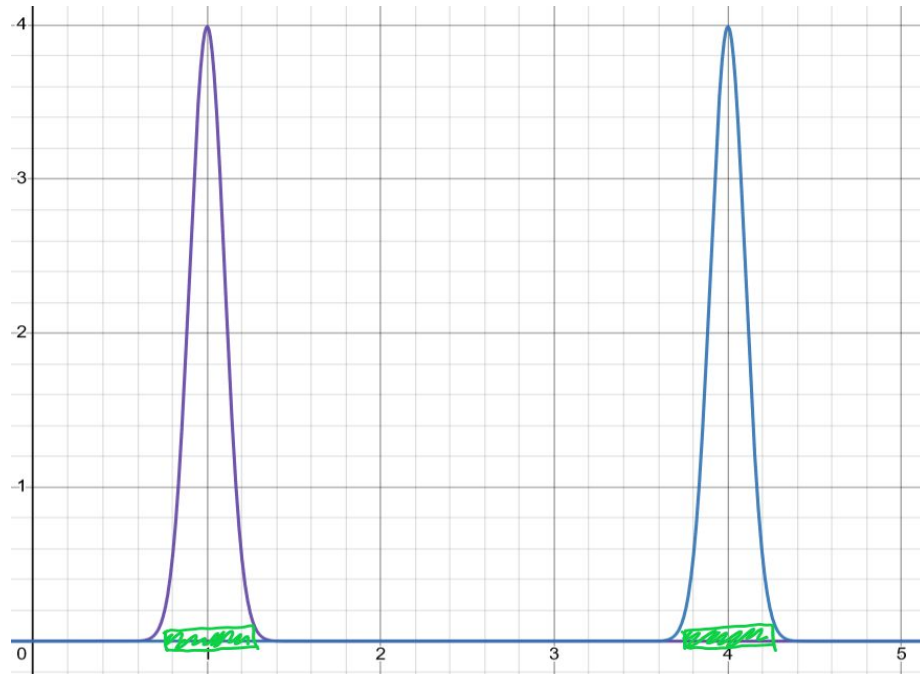
- Probabilistic **programs** define p
- Running the program gives the **prior** $p(\mathbf{x})$
- We are interested in the **posterior** $p(\mathbf{y} \mid \mathbf{x})$
- Estimating **marginal likelihood** $p(\mathbf{x} \mid \mathbf{y})$ precisely is difficult

Why is sampling $p(x | y)$ difficult?



Source: https://rpubs.com/hillt5/blog2_621

Why is sampling $p(x | y)$ difficult?



b)

Introduction

- What if we could "build" / "approximate" $p(x | y)$ in some other way?

Mean Field Approximation

- The probability of a trace is given by:

$$p(x) = \prod_{t=1}^T p_t(x_t \mid \psi_t(h_t))$$

where h_t is the history (x_1, \dots, x_{t-1}) of the program up to ERP t

Mean Field Approximation

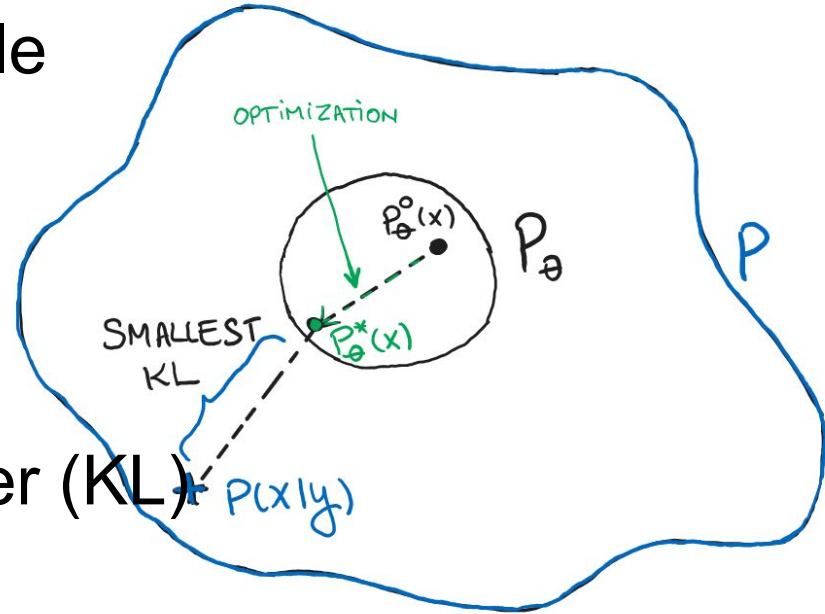
- We approximate $p(x)$ with a simpler $p_\theta(x)$
- We use θ to parameterize p_θ , then we learn θ

$$p(x) = \prod_{t=1}^T p_t(x_t | \psi_t(h_t))$$

$$p_\theta(x) = \prod_{t=1}^T p_\theta(x_t | \theta_t)$$

Mean Field Approximation

- Analytically intractable
- Stochastic Gradient Descent is our friend
- Enter Kullback-Leibler (KL) Divergence!



KL Divergence

- Measure of distance between distributions
- Essentially our reward function

$$KL(p_{\theta}, p(x|y)) = \int_x p_{\theta}(x) \log \left(\frac{p_{\theta}(x)}{p(x|y)} \right)$$

Bayes' rule

$$= \int_x p_{\theta}(x) \log \left(\frac{p_{\theta}(x)}{p(y|x)p(x)} \right) + \log p(y) =$$

$$= \underbrace{-L(\theta)}_{\text{ELBO}} + \log p(y)$$

KL Divergence and ELBO

$$KL(p_\theta, p(x|y)) = -L(\theta) + \log p(y)$$

where

$$L(\theta) \triangleq \int_x p_\theta(x) \log \left(\frac{p(y|x)p(x)}{p_\theta(x)} \right)$$

- $KL \geq 0 \rightarrow -L(\theta) + \log(p(y)) \geq 0$
- $\log(p(y)) \geq L(\theta)$
(constant)

Stochastic Gradient Optimization

$$-\nabla_{\theta} L(\theta) = \int_x \nabla_{\theta} \left(p_{\theta}(x) \log \left(\frac{p_{\theta}(x)}{p(y|x)p(x)} \right) \right) \quad (5)$$

$$\stackrel{\text{product rule}}{=} \int_x \nabla_{\theta} p_{\theta}(x) \left(\log \left(\frac{p_{\theta}(x)}{p(y|x)p(x)} \right) \right) + \int_x p_{\theta}(x) (\nabla_{\theta} \log(p_{\theta}(x))) \quad (6)$$

Stochastic Gradient Optimization

- To derive the MC estimate, we use a few tricks

$$\text{a) } \nabla \log p_{\theta}(x) = \frac{\nabla_{\theta} p_{\theta}(x)}{p_{\theta}(x)}$$

$$\text{b) } \int_x p_{\theta}(x) \nabla \log p_{\theta}(x) = 0$$

Stochastic Gradient Optimization

$$-\nabla_{\theta} L(\theta) = \int_x \nabla_{\theta} p_{\theta}(x) \left(\log \left(\frac{p_{\theta}(x)}{p(y|x)p(x)} \right) \right) + \int_x p_{\theta}(x) (\nabla_{\theta} \log(p_{\theta}(x))) \quad (6)$$

$$= \int_x \nabla_{\theta} p_{\theta}(x) \left(\log \left(\frac{p_{\theta}(x)}{p(y|x)p(x)} \right) \right) \quad (7)$$

$$\stackrel{\text{apply a)}}{=} \int_x p_{\theta}(x) \nabla_{\theta} \log(p_{\theta}(x)) \left(\log \left(\frac{p_{\theta}(x)}{p(y|x)p(x)} \right) \right) \quad (8)$$

$$= \int_x p_{\theta}(x) \nabla_{\theta} \log(p_{\theta}(x)) \left(\log \left(\frac{p_{\theta}(x)}{p(y|x)p(x)} \right) + K \right) \quad (9)$$

$$\approx \frac{1}{N} \sum_{x^j} \nabla_{\theta} \log p_{\theta}(x^j) \left(\log \left(\frac{p_{\theta}(x^j)}{p(y|x^j)p(x^j)} \right) + K \right) \quad (10)$$

$$\text{a) } \nabla \log p_{\theta}(x) = \frac{\nabla_{\theta} p_{\theta}(x)}{p_{\theta}(x)}$$

Can now use MC estimate

$$\text{b) } \int_x p_{\theta}(x) \nabla \log p_{\theta}(x) = 0$$

Mean-field approximation

Probabilistic program A

```
1: M = normal();
2: if M>1
3:   mu = complex_deterministic_func( M );
4:   X = normal( mu );
5: else
6:   X = rand();
7: end;
```

Mean-Field variational program A

```
1: M = normal(  $\theta_1$  );
2: if M>1
3:   mu = complex_deterministic_func( M );
4:   X = normal(  $\theta_3$  );
5: else
6:   X = rand( $\theta_4, \theta_5$ );
7: end;
```

Compositional Variational Inference

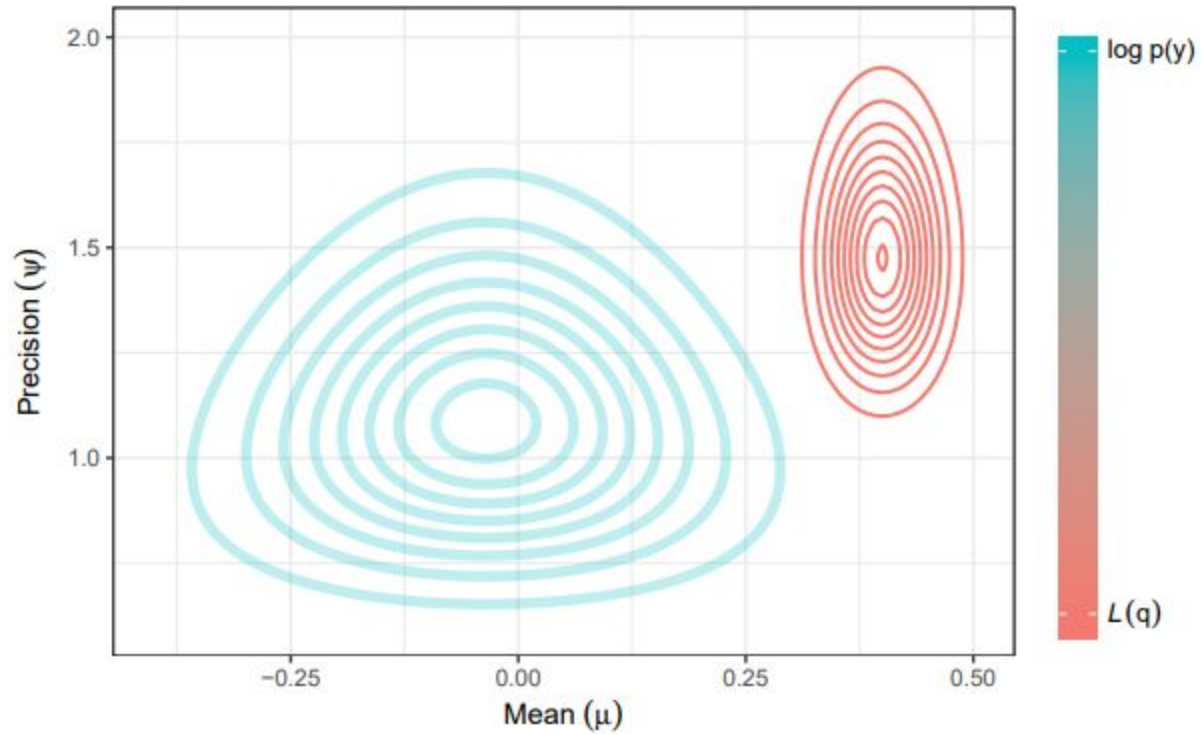
- Initialize θ to arbitrary value
- Sample $x_t \sim p_\theta(x_t)$
- Compute:
 - $\log(p_\theta(x_t))$
 - $\log(p(x_t | h_t))$
 - $R_t = \log(p(x_t | h_t)) - \log(p_\theta(x_t))$
 - local gradient $\Psi_t = \nabla_{\theta_t} \log(p_{\theta_t}(x_t))$

Computing the gain

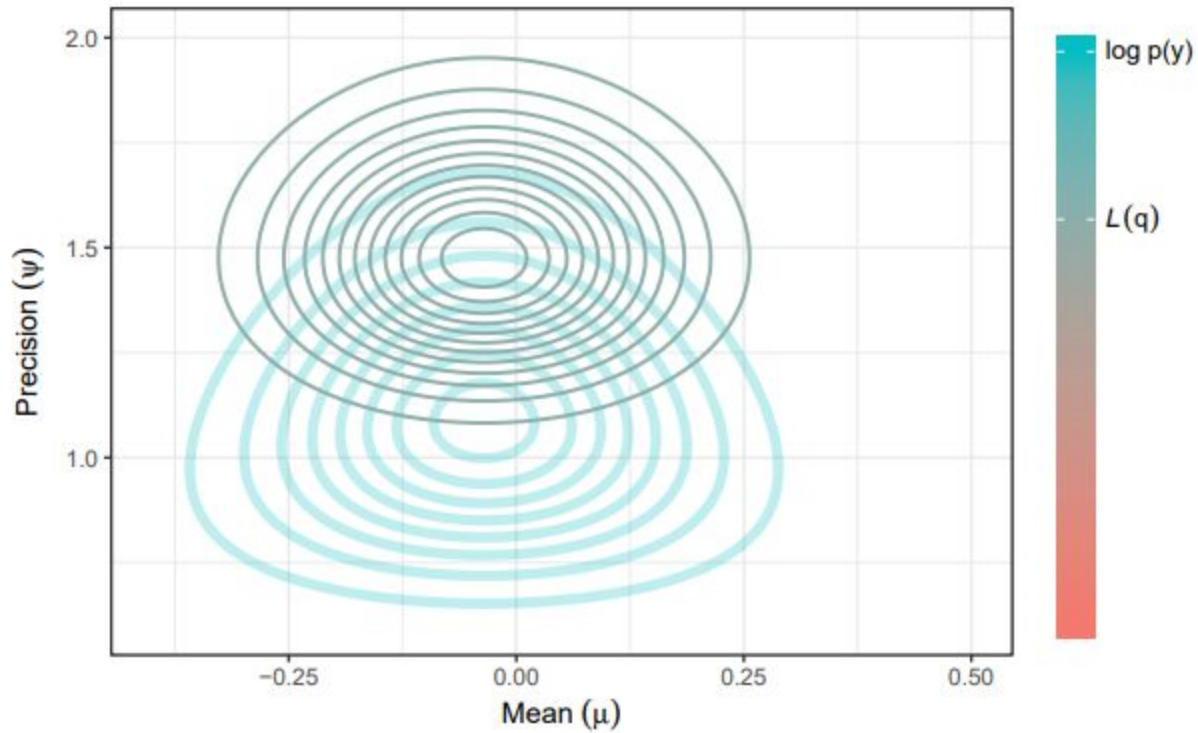
- Compute
 - $\log p(y | x)$
 - the gain $R = \sum R_t + \log p(y | x) + K$
- Estimate at ERP t can be averaged over many sample traces for a more accurate estimate

$$-\nabla_{\theta} L(\theta) \approx \frac{1}{N} \sum_{x^j} \nabla_{\theta} \log p_{\theta}(x^j) \left(\log \left(\frac{p_{\theta}(x^j)}{p(y|x^j)p(x^j)} \right) + K \right) \quad (10)$$

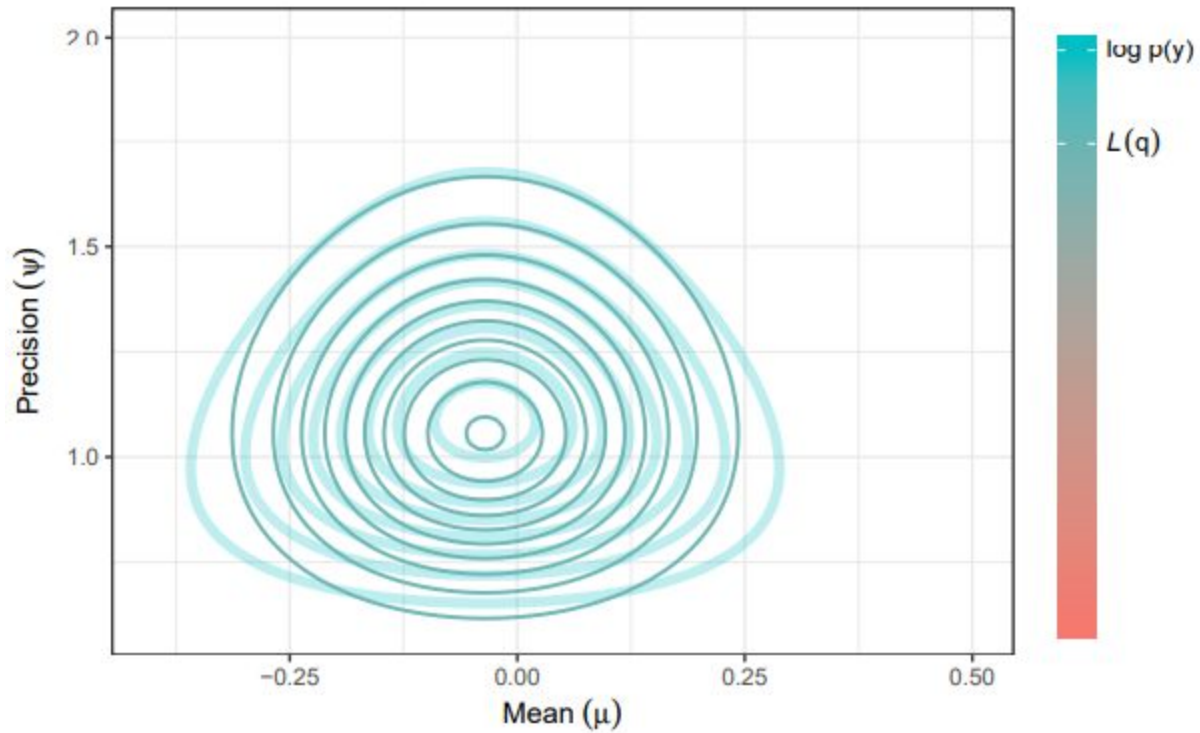
Iteration 0 (initialisation)



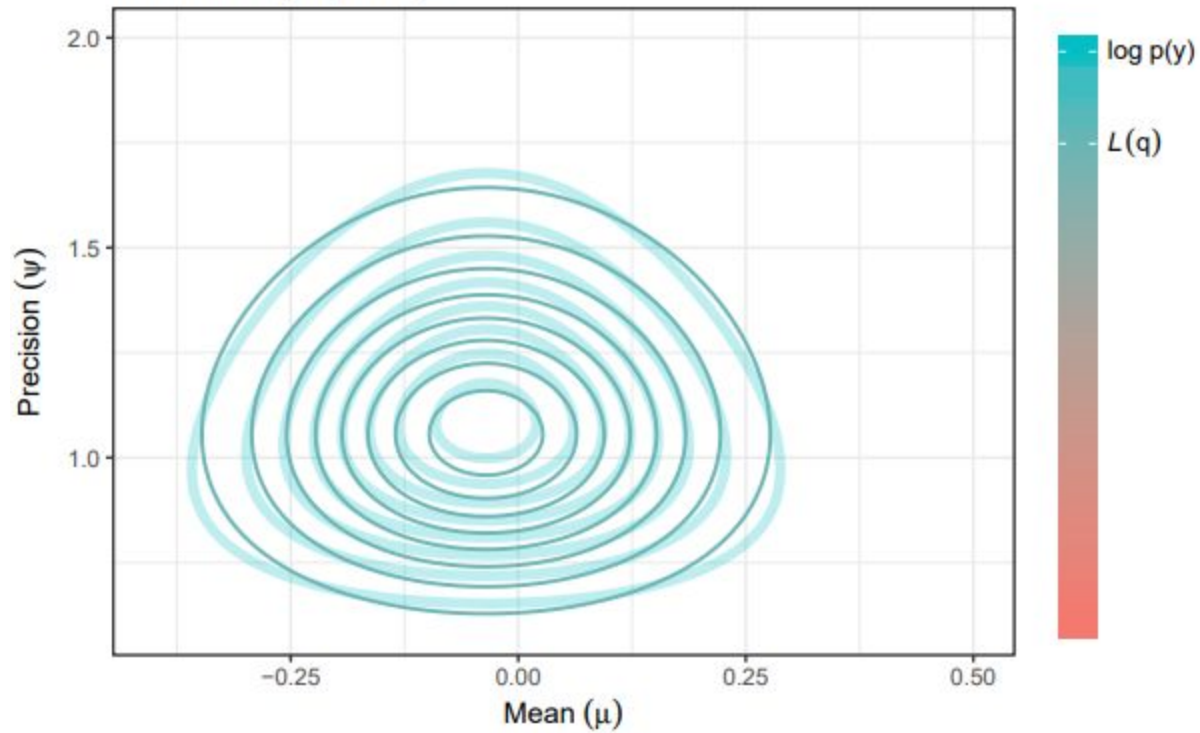
Iteration 1 (μ update)



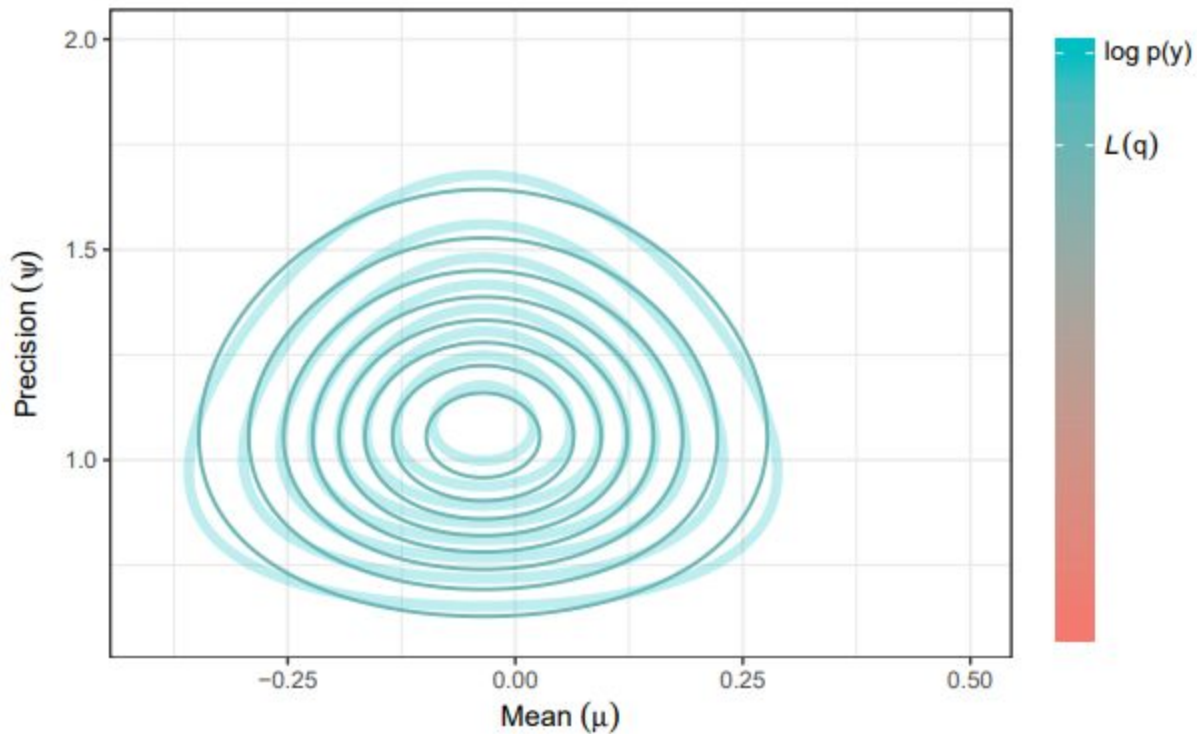
Iteration 1 (ψ update)



Iteration 2 (μ update)



Iteration 2 (ψ update)

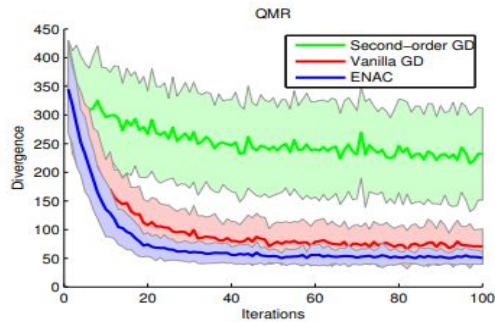


Experiments

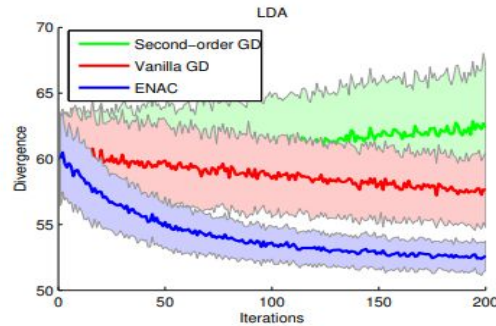
- Which algorithm estimates the direction of the gradient best? How does SGD compare with the others?

Experiments

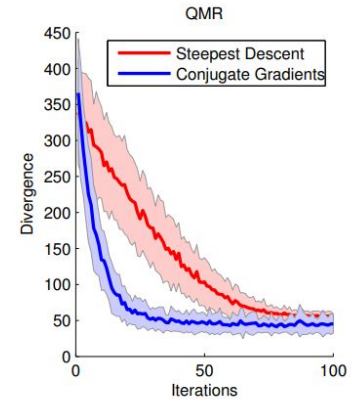
- ENAC > Vanilla GD (ours) > SOGD



(a) Results on QMR



(b) Results on LDA



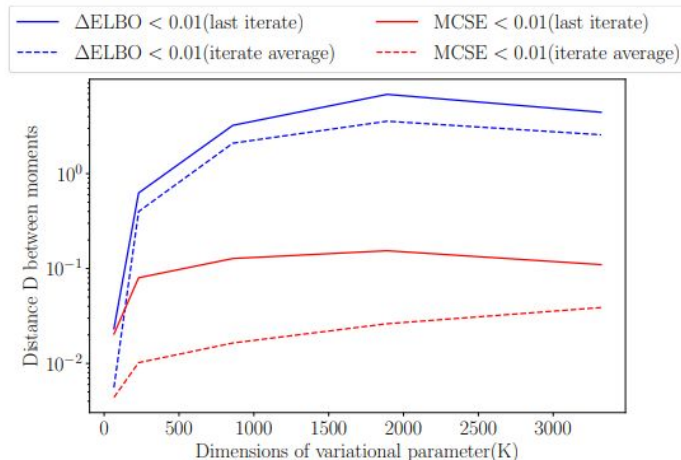
(c) Steepest descent vs. conjugate gradients

Takeaways

- Very fast approximate sampling from the posterior
- Much cheaper than using MCMC sampling

Related Work

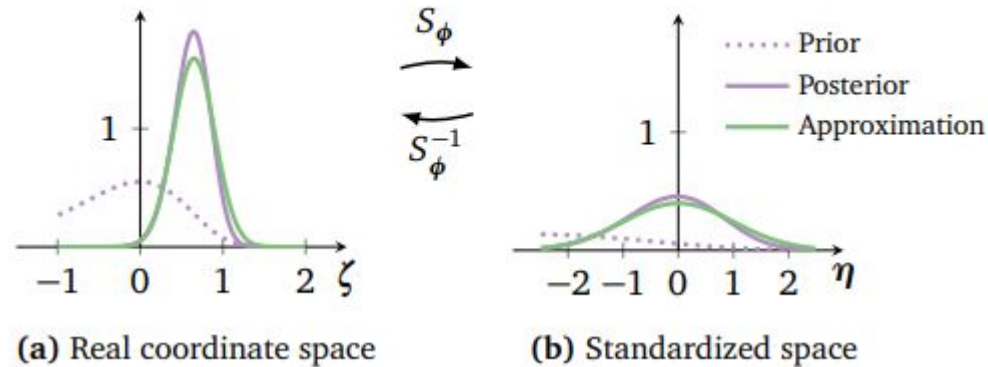
- High-dimensional posteriors are sometimes poorly approximated by SGD, a parallelizable approach is proposed [1]



$$\bar{\lambda} \equiv \frac{1}{T} \sum_{i=1}^T \lambda_{t+i},$$

Related Work

- Automatic Differentiation Variational Inference [2]



Related Work

- Rethinking Variational Inference for Probabilistic Programs with Stochastic Support [3]

Probabilistic program A

```
1: M = normal();
2: if M>1
3:   mu = complex_deterministic_func( M );
4:   X = normal( mu );
5: else
6:   X = rand();
7: end;
```

Questions?

References

- [1] - Dhaka, A. K., Catalina, A., Andersen, M. R., Magnusson, M., Huggins, J., & Vehtari, A. (2020). Robust, accurate stochastic optimization for variational inference. *Advances in Neural Information Processing Systems*, 33, 10961-10973.
- [2] - Ambrogioni, L., Lin, K., Fertig, E., Vikram, S., Hinne, M., Moore, D., & van Gerven, M. (2021, March). Automatic structured variational inference. In *International Conference on Artificial Intelligence and Statistics* (pp. 676-684). PMLR.
- [3] - Reichelt, T., Ong, L., & Rainforth, T. (2022). Rethinking Variational Inference for Probabilistic Programs with Stochastic Support. *Advances in Neural Information Processing Systems*, 35, 15160-15175.