

STOCHASTIC VARIATIONAL INFERENCE FOR LARGE SCALE MACHINE LEARNING.

①

VARIATIONAL BAYES: APPROXIMATE POSTERIOR

$$P(\tilde{z}, \tilde{\theta} | \tilde{x}) = \frac{\prod_{i=1}^M P(z_i, \theta_i | \tilde{x}_i) P(\tilde{\theta})}{P(\tilde{x})}$$

WITH SIMPLER DISTRIBUTION $q(\tilde{z}, \tilde{\theta})$. PARAMETERIZED BY
~~A~~ A SET OF PARAMETERS λ .

VB MAXIMIZES THE EVIDENCE LOWER BOUND OR ELBO

$$\begin{aligned} \log P(\tilde{x}) &= \log \mathbb{E}_{q(\lambda)} \left[\frac{P(\tilde{x}, \tilde{z}, \tilde{\theta})}{q(\tilde{z}, \tilde{\theta})} \right] \geq \\ &\geq \underbrace{\mathbb{E}_{q(\lambda)} \left[\log P(\tilde{x}, \tilde{z}, \tilde{\theta}) \right] - \mathbb{E}_{q(\lambda)} \left[\log q(\tilde{z}, \tilde{\theta}) \right]}_{\mathcal{L}(\lambda)} \end{aligned}$$

LET'S ASSUME THAT

$$\mathcal{L}(\lambda) = \frac{1}{M} \sum_{i=1}^M f(\tilde{x}_i, \lambda_i)$$

AND THAT λ IS A VECTOR. WHAT TECHNIQUES
CAN WE USE TO OPTIMIZE $\mathcal{L}(\lambda)$?

GRADIENT DESCENT (GD)

EACH ITERATION UPDATES λ USING

$$\lambda_{t+1} = \lambda_t + \gamma \frac{1}{M} \sum_{i=1}^M \nabla_{\lambda} f(\tilde{x}_i, \lambda_t)$$

UNDER SUFFICIENT REGULARITY ASSUMPTIONS, ~~GD~~
WHEN γ IS SMALL ENOUGH, GD HAS
LINEAR CONVERGENCE: $\|y_t - y^*\| \sim \rho^t$ WHERE ρ
IS THE RESIDUAL ERROR: $\rho = |\lambda_t - \lambda^*|$

SECOND ORDER GRADIENT DESCENT (ZGD)

EACH ITERATION UPDATES λ USING

$$\lambda_{t+1} = \lambda_t + \eta_t \frac{1}{M} \sum_{i=1}^M \nabla_{\lambda} f(\tilde{x}_i, \lambda_t)$$

APPROACHES
THE INVERSE
OF THE
HESIAN
AT THE
OPTIMUM.

UNDER SUFF. REGULARITY ASSUMPTIONS ZGD HAS
QUADRATIC CONVERGENCE: $\|y_t - y^*\| \sim \rho^2$

STOCHASTIC GRADIENT DESCENT.

$$\lambda_{t+1} = \lambda_t + \gamma_t \nabla_{\lambda} f(\tilde{x}_{z_t}, \lambda_t)$$

WHERE z_t IS SAMPLED UNIFORMLY IN $\{1, \dots, M\}$.

CONVERGENCE REQUIRES

$$\sum_t \gamma_t^2 < \infty$$

$$\sum_t \gamma_t = \infty$$

WE CAN CHOOSE $\delta \in (0, 5, 1]$ UNDER SOME CONDITIONS THE BEST CONVERGENCE SPEED IS ACHIEVED WHEN $\delta = (\frac{1}{\lambda} \frac{1}{\sigma})^{-2}$ WITH REGULARITY

$\delta_2 \sim t^{-1}$ IN THIS CASE $E[p] \sim t^{-1}$.

ASYMPTOTIC BEHAVIOR OF EACH METHOD:

	GD	ZGD	SGD
TIME PER ITERATION	M	M	1
ITERATIONS TO ACCURACY ρ	$\log \frac{1}{\rho}$	$\log \log \frac{1}{\rho}$	$\frac{1}{\rho}$
TIME TO ACCURACY ρ	$M \log \frac{1}{\rho}$	$M \log \log \frac{1}{\rho}$	$\frac{1}{\rho}$

↓ DOES NOT DEPEND ON M!!

LET'S APPLY SGD TO OUR VARIATIONAL PROBLEM.

WE ASSUME CONJUGATE EXPONENTIAL FAMILIES FOR OUR PRIOR AND LIKELIHOOD.

$$P(\theta) \propto \exp \left\{ \eta_{\theta}^T g_{\theta} \right\}$$

$$P(x_i, z_i | \theta) \propto \exp \left\{ \eta_{\theta}^T t(x_i, z_i) \right\}$$

THE CONDITIONALS ARE IN THE SAME EXPONENTIAL FAMILY AS THE PRIOR.

$$P(z_i | x_i, \theta) \propto \exp \left\{ \eta_{\theta}^T t(x_i, z_i) \right\}$$

$$P(\theta | x, z) \propto \exp \left\{ \eta_{\theta}^T \left(\alpha + \sum_{i=1}^n t(x_i, z_i) \right) \right\}$$

$$q_{\lambda}(\theta, z) = q_{\lambda}(\theta) \prod q_{\lambda}(z_i)$$

WE CHOOSE ~~FOR~~ q_{λ} ~~THE~~ A FACTORIZED FORM (MEAN-FIELD). THE OPTIMAL FACTORS FOR q_{λ} ARE IN THE SAME FAMILY AS THE COMPLETE CONDITIONALS WHY?

$$q(z_i) = \exp \left\{ \beta_i^T t(z_i) - z(\beta_i) \right\}$$

$$q(\theta | \lambda) = \exp \left\{ \eta_{\theta}^T \lambda - z(\lambda) \right\}$$

$$\lambda = \{ \eta_{\theta}, \beta_1, \dots, \beta_n \}$$

AS A FUNCTION OF λ THE ELBO IS

$$\mathcal{L}(\lambda) = \mathbb{E}_q \left[\log p(\theta, x, z) \right] - \mathbb{E}_q \left[\log q(\theta) \right] + \text{const.}$$

$$\mathcal{L}(\theta) = \mathbb{E}_g \left[\alpha + \sum_{i=1}^M \frac{t}{n} (x_i, z_i) \right] \left[\nabla_{\theta} z(\theta) \right] - \theta^T \nabla_{\theta} z(\theta) + z(\theta) + \text{CONST.}$$

$$\nabla_{\theta} z(\theta) = \mathbb{E}_g [g(\theta)]$$

AND ITS GRADIENT IS

$$\nabla_{\theta} \mathcal{L}(\theta) = \nabla_{\theta}^2 z(\theta) \left\{ \mathbb{E}_g \left[\alpha + \sum_{i=1}^M \frac{t}{n} (x_i, z_i) \right] - \theta \right\}$$

THE ELBO IS OPTIMIZED W.R.T. θ BY

SETTING $\theta = \mathbb{E}_g \left[\alpha + \sum_{i=1}^M \frac{t}{n} (x_i, z_i) \right]$

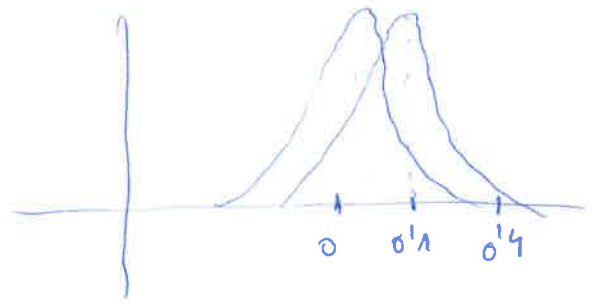
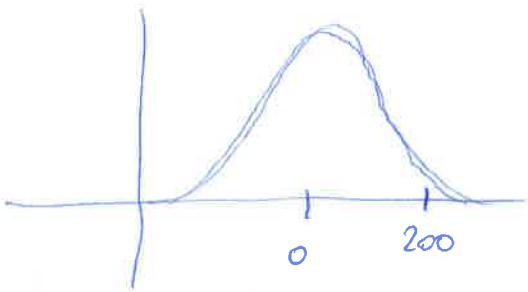
W.R.T. β_i BY SETTING $\beta_i = \mathbb{E}_g \left[\frac{t}{n} (x_i, z_i) \right]$

SLOW CONVERGENCE !!

ALTERNATIVE : SGD. USE NATURAL GRADIENTS INSTEAD OF STANDARD GRADIENTS.

EUCLIDIAN DISTANCE BETWEEN $N(0, 10^4)$ AND $N(10, 10^4)$ IS 10.

DISTANCE BETWEEN $N(0, 10^{-2})$ AND $N(10^{-1}, 10^{-2})$ IS 10^{-1}



EUCLIDIAN DISTANCE BETWEEN $q(\theta | \Sigma)$ AND $q(\theta | \Sigma')$
 IS $\|\theta - \theta'\|_1$

USE KL-DISTANCE: ~~KL-DISTANCE~~

↓
 SYMMETRIZED DIST $[q(\theta | \Sigma), q(\theta | \Sigma')] =$

$$= \mathbb{E}_{q_{\Sigma}} \left[\log \frac{q(\theta | \Sigma)}{q(\theta | \Sigma')} \right] + \mathbb{E}_{q_{\Sigma'}} \left[\log \frac{q(\theta | \Sigma')}{q(\theta | \Sigma)} \right]$$

THE DIRECTION OF THE NATURAL GRADIENT IS

ARG MAX. $\mathcal{L}(\theta + \Delta\theta)$ s.t. $\text{KL}^{\text{SYM}}(\theta, \theta + \Delta\theta) < \epsilon$

$\Delta\theta$

WHEN $\epsilon \rightarrow 0$.

THE NATURAL GRADIENT CAN BE SHOWN TO BE

$$\nabla_{\theta}^* \mathcal{L}(\theta) = G(\theta)^{-1} \nabla_{\theta} \mathcal{L}(\theta)$$

RIEMANNIAN METRIC TENSOR

$$\Delta\theta^T G(\theta) \Delta\theta =$$

$$\begin{aligned}
\text{KL}^{\text{SYM}}(\tilde{x}, \tilde{x} + \Delta\tilde{x}) &= \cancel{E_{\tilde{x}} \left[\frac{z(\tilde{x})}{z(\tilde{x} + \Delta\tilde{x})} \right]} \\
&= \cancel{\nabla z(\tilde{x})^T \tilde{x}} - \cancel{z(\tilde{x})} - \cancel{\nabla z(\tilde{x})^T (\tilde{x} + \Delta\tilde{x})} + \cancel{z(\tilde{x} + \Delta\tilde{x})} \\
&\quad + \cancel{\nabla z(\tilde{x} + \Delta\tilde{x})^T (\tilde{x} + \Delta\tilde{x})} - \cancel{z(\tilde{x} + \Delta\tilde{x})} - \\
&\quad - \cancel{\nabla z(\tilde{x} + \Delta\tilde{x})^T \tilde{x}} + \cancel{z(\tilde{x})} \\
&= -\nabla z(\tilde{x})^T \Delta\tilde{x} + \nabla z(\tilde{x} + \Delta\tilde{x})^T \Delta\tilde{x} \\
&= \left[\nabla z(\tilde{x} + \Delta\tilde{x})^T - \nabla z(\tilde{x})^T \right] \Delta\tilde{x} \\
&= \cancel{\Delta\tilde{x}^T \nabla z(\tilde{x}) \Delta\tilde{x}} \\
&= \cancel{\Delta\tilde{x}^T \left[\nabla z(\tilde{x}) + \Delta\tilde{x} \nabla^2 z(\tilde{x}) \right] \Delta\tilde{x}} - \cancel{\nabla z(\tilde{x})^T} \Delta\tilde{x} = \\
&= \Delta\tilde{x}^T \nabla^2 z(\tilde{x}) \Delta\tilde{x}
\end{aligned}$$

THEREFORE $G_{\tilde{x}}^{\text{NOISY}}(\tilde{x}) = \nabla^2 z(\tilde{x})$.

AND $\nabla_{\tilde{x}}^* \mathcal{L}(\tilde{x}) = E_{\tilde{g}} \left[\tilde{x} + \sum_{i=1}^M \frac{t(\tilde{x}_i, z_i)}{\tilde{x}} \right] - \tilde{x}$

NOISY ESTIMATE OF THE NATURAL GRADIENT.

$$\nabla_{\tilde{x}}^* \text{NOISY} \mathcal{L}(\tilde{x}) = \tilde{x} + \frac{1}{P(\tilde{x})} E_{\tilde{g}} \left[\frac{t(\tilde{x}_i, z_i)}{\tilde{x}} \right] - \tilde{x}$$

WHERE $p(k)$ IS THE PROB. OF SUBSAMPLING THE k -TH DATA POINT.

THE SVI UPDATE FOR $\hat{\theta}$ IS

$$\begin{aligned}\hat{\theta}_k^{\text{new}} &= \hat{\theta}_k^{\text{old}} + p_k \nabla_{\hat{\theta}}^* \text{NOISY} \ell(\hat{\theta}) = \\ &= (1 - p_k) \hat{\theta}_k^{\text{old}} + p_k \hat{\theta}_k^{\text{noisy}}\end{aligned}$$

WHERE $\hat{\theta}_k^{\text{noisy}} = \hat{\theta}_k + \frac{1}{p(k)} \mathbb{E}_{\mathcal{G}} \left[\ell_{\hat{\theta}}(x_k, z_k) \right]$

OPTIMAL PARAMETER IF THE DATA INCLUDES ONLY DATA k ~~WITH PROB. $p(k)$~~ REPLICATED $\frac{1}{p(k)}$ TIMES.

WE CAN USE MINIBATCHES OF SIZE S TO REDUCE VARIANCE. IN THIS CASE

$$\hat{\theta}_k^{\text{new}} = (1 - p_k) \hat{\theta}_k^{\text{old}} + \frac{1}{S} \sum_{i=1}^S \hat{\theta}_k^{\text{noisy}}(i)$$

SVI PSEUDO CODE:

REPEAT

SAMPLE x_k WITH PROB. $p(k)$.

OPTIMALLY UPDATE $\hat{\theta}_k = \mathbb{E}_{\mathcal{G}} \left[\nabla_{\hat{\theta}} \ell_{\hat{\theta}}(x_k, z_k) \right]$

COMPUTE NOISY OPTIMAL $\hat{\theta}$

$$\hat{\theta}_k^{\text{noisy}} = \hat{\theta}_k + \frac{1}{p(k)} \mathbb{E}_{\mathcal{G}} \left[\ell_{\hat{\theta}}(x_k, z_k) \right]$$

UPDATE \tilde{z} USING $\tilde{z}^{new} = (1 - \rho_t) \tilde{z}^{old} + \rho_t \tilde{z}^{noisy}$

UNTIL FOREVER.

$$\rho_t = (t + \tau)^{-\kappa} \quad \kappa \in (0's, 1] , \tau \geq 0.$$

SVI IN PROB. BAYESIAN MATRIX FACTORIZATION.

$$\boxed{M} \approx \boxed{U} \times \boxed{V^T}$$

M IS AN $m \times d$ MATRIX

U AND V ARE $m \times k$ AND $d \times k$ MATRICES

$$P(M | U, V) = \prod_{i=1}^m \prod_{j=1}^d N(m_{ij} | u_i \cdot v_j^T, \sigma^2)$$

$$\left[\prod_{i=1}^m \prod_{l=1}^k N(u_{i,l} | 0, \tau^2) \right] \left[\prod_{j=1}^d \prod_{l=1}^k N(v_{j,l} | 0, \rho^2) \right]$$

THE COMPLETE CONDITIONALS ARE GAUSSIANS.
THE LIKELIHOOD AND THE PRIOR ARE ~~CONJUGATE~~ CONJUGATE.

$$Q(U, V) = \left[\prod_{i=1}^m N(u_i | m_i^U, \tau_i^U) \right]$$

$$\left[\prod_{j=1}^d N(v_j | m_j^V, \tau_j^V) \right]$$

WE WORK WITH THE NATURAL PARAMETERS

$$h_{z,i}^u = \left[\begin{matrix} v_i^u \\ \tilde{z} \end{matrix} \right]^{-1} \quad \text{AND} \quad \cancel{h_{z,i}^u} = \left[\begin{matrix} v_i^u \\ \tilde{z} \end{matrix} \right]^{-1} m_i^u$$

THE BATCH UPDATES ARE

$$h_{z,i}^u = \left[\begin{matrix} d \\ \sum_{j=1}^d \cancel{v_j^u} \end{matrix} \left[\begin{matrix} v_j^u \\ \tilde{z} \end{matrix} \right] + \frac{m_j^u m_j^{uT}}{\sigma^2} \right] + 1/\sigma^2 I_{\tilde{z}}$$

$$h_{\lambda,i}^u = \sum_{j=1}^d \frac{m_{ij}^u m_j^u}{\sigma^2}$$

THE SVI UPDATES ARE:

$$h_{z,i}^u \text{ new} = h_{z,i}^u \text{ old} (1 - \rho_t) + \rho_t h_{z,i}^u \text{ NOISY MINIBATCH}$$

$$h_{\lambda,i}^u \text{ new} = h_{\lambda,i}^u \text{ old} (1 - \rho_t) + \rho_t h_{\lambda,i}^u \text{ NOISY MINIBATCH}$$

WHERE

$$h_{z,i}^u \text{ NOISY MINIBATCH} = \frac{d}{|I|} \left[\sum_{j \in I} \frac{v_j^u}{\sigma^2} + \frac{m_j^u m_j^{uT}}{\sigma^2} \right] + \frac{1}{\sigma^2} I_{\tilde{z}}$$

$$h_{\lambda,i}^u \text{ NOISY MINIBATCH} = \frac{d}{|I|} \left[\sum_{j \in I} \frac{m_{ij}^u m_j^u}{\sigma^2} \right]$$

RANDOM

WHERE I IS A SET OF ~~RANDOM~~ COLUMN INDEXES, SUCH THAT ~~WE CAN~~ ~~WRITE~~ ~~THE~~ ~~PROBABILITIES~~ EACH ELEMENT OF I IS SAMPLED UNIFORMLY WITH REPLACEMENT FROM $\{1, \dots, d\}$.

USEFULL REFERENCES:

VARIATIONAL BATES:

- 1999 ATTIAS UAI.
- 2001 GHAMRANI, BEAL. ADVANCE MEAN FIELD METHOD — THEORY AND PRACTICE, CHAPTER GRAPHICAL MODELS AND VARIATIONAL METHODS.
- 2006 BISHOP PATTERN RECOGNITION AND MACHINE LEARNING -

STOCHASTIC GRADIENT DESCENT:

- 2010 BOTTOU LARGE SCALE MACHINE LEARNING WITH STOCHASTIC GRADIENT DESCENT.
- 1998 BOTTOU ON LINE, LEARNING AND STOCHASTIC APPROXIMATIONS
- 1983 DENNIS SCHNABEL NUMERICAL METHODS FOR UNCONSTRAINED OPTIMIZATION AND NON-LINEAR EQUATIONS.

STOCHASTIC VARIATIONAL INFERENCE

2013 HOFFMAN, BLEI, CHONG, PAISLEY
STOCHASTIC VARIATIONAL INFERENCE

NATURAL GRADIENTS

1998 AMARI
NATURAL GRADIENTS WORKS
EFFICIENTLY IN LEARNING