

title of the big homework

Big homework in big course

Saurav Shekhar, (16-947-921)

shekhars@student.ethz.ch

November 14, 2018

Abstract

All notes regarding my masters thesis project

1 External links

- Thesis pre proposal
- Tasks document
- Gitlab repository might be removed
- Overleaf project not regularly updated

2 Math notes

Taylor series multivariate.

$$f(\mathbf{r}_0 + \mathbf{a}t) = f(\mathbf{r}_0) + [\mathbf{a} \cdot \nabla f(\mathbf{r})] \Big|_{\mathbf{r}=\mathbf{r}_0} t + \frac{1}{2!} [\mathbf{a} \cdot \nabla][\mathbf{a} \cdot \nabla] f(\mathbf{r}) \Big|_{\mathbf{r}=\mathbf{r}_0} t^2 + \dots \quad (1)$$

Shannon Entropy.

Self information of event $x = x$ is defined as $I(x) := -\log P(x)$

$$H(x) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\log P(x)]$$

Cramer-Rao lower bound. [Balabdaoui and van de Geer, 2016] Suppose θ is an unknown deterministic parameter which is to be estimated from measurements x , distributed according to some pdf $f(x; \theta)$. The variance of any unbiased estimator $\hat{\theta}$ of θ is then bounded by reciprocal of Fischer Information $I(\theta)$:

$$\begin{aligned} \text{var}(\hat{\theta}) &\geq \frac{1}{I(\theta)} \text{ where} \\ I(\theta) &= \mathbb{E} \left[\left(\frac{\partial l(x; \theta)}{\partial \theta} \right)^2 \right] \\ &= -\mathbb{E} \left[\frac{\partial^2 l(x; \theta)}{\partial \theta^2} \right] \end{aligned}$$

Note: See Wikipedia for other more general versions

Lipschitz continuity, smoothness, etc.. f is M -Lipschitz continuous given M if

$$|f(x) - f(y)| \leq M|x - y| \forall x, y \in \mathbb{R}$$

. If f is differentiable, Lipschitz continuity says that f has bounded derivative.

f is L -Lipschitz smooth if its derivatives are Lipschitz continuous with L . This is called smoothness type $C^{1,1}$ i.e

$$\forall x, y \in \mathbb{R}, \|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$$

The definition does not assume convexity of f . Some other equivalent conditions are: [on here](#).

$$g(x) = \frac{L}{2}x^T x - f(x) \text{ is convex} \quad (2)$$

$$f(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{L}{2}\|y - x\|^2, \forall x, y \quad (3)$$

$$(\nabla f(x) - \nabla f(y))^T(x - y) \leq L\|x - y\|^2, \forall x, y \quad (4)$$

$$f(\alpha x + (1 - \alpha)y) \geq \alpha f(x) + (1 - \alpha)f(y) - \frac{\alpha(1 - \alpha)L}{2}\|x - y\|^2 \forall x, y \in \mathbb{R}, \alpha \in [0, 1] \quad (5)$$

$$\dots \quad (6)$$

Strong Convexity: f is α -strongly convex if

$$\nabla^2 f(x) \succeq \alpha I \forall x$$

Also

$$f(x + y) \geq f(x) + y^T \nabla f(x) + \frac{\alpha}{2}\|x - y\|^2 \quad (7)$$

7 is equivalent to saying $g(x) = f(x) - \frac{\alpha}{2}\|x\|^2$ is convex.. Latter is equivalent to $\nabla^2 g \succeq 0 \equiv \nabla^2 f \succeq \alpha I$. [more info here](#)

Note:- strong convexity can be written as

$$f(y) \geq f(x) + \nabla f(x)^T(y - x) + \frac{l}{2}\|y - x\|^2 \quad (8)$$

Equation 8 is direct contrast to 3. $\frac{L}{l}$ is called condition number of matrix.

3 Code Notes

Normal distribution edward

```
1  from edward.models import Normal
2  from keras.layers import Dense
3
4  hidden = Dense(256, activation='relu')(x_ph)
5  qz = Normal(loc=Dense(10)(hidden),
6  scale=Dense(10, activation='softplus')(hidden))
```

4 Web pages

1. [Zen of gradient descent with intro to Nesterov Method](#)
2. [I am a bandit Nesterov accelerated](#)
3. [CMU Stats FW lecture](#)

5 Literature notes

Variational inference.: [Blei et al., 2016] **INPROGRESS**

$$\underbrace{\text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}))}_{\downarrow \text{Objective} \geq 0} := \mathbb{E}_{q(\mathbf{z})} [\log q(\mathbf{z})] - \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{z}|\mathbf{x})] \quad (9)$$

$$\text{ELBO}(q) := \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_{q(\mathbf{z})} [\log q(\mathbf{z})] \quad (10)$$

$$\text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) + \underbrace{\text{ELBO}(q)}_{\uparrow \text{Optimize}} = \underbrace{\log p(\mathbf{x})}_{\text{constant w. } q} \quad (11)$$

$$\text{ELBO}(q) = \mathbb{E} [\log p(\mathbf{x}|\mathbf{z})] - \text{KL}(q(\mathbf{z})||p(\mathbf{z})) \quad (12)$$

$$\text{ELBO}(q) = \mathcal{Q}(\theta, \theta_t) - \mathcal{H}(\mathbf{z}|\mathbf{x}) \text{ //Entropy} \quad (13)$$

$$(14)$$

Algorithm 1: Coordinate Ascent for VI

Input: A model $p(\mathbf{x}, \mathbf{z})$, a data set \mathbf{x}

Output: A variational density $q(\mathbf{z}) = \prod_{j=1}^m q_j(z_j)$

```

1 Initialize: Variational factors  $q_j(z_j)$ 
2 while the ELBO has not converged do
3   for  $j \in \{1, \dots, m\}$  do
4     Set  $q_j(z_j) \propto \exp\{\mathbb{E}_{-j}[\log p(z_j | \mathbf{z}_{-j}, \mathbf{x})]\}$ 
5   end
6   Compute  $\text{ELBO}(q) = \mathbb{E} [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E} [\log q(\mathbf{z})]$ 
7 end
8 return  $q(\mathbf{z})$ 
```

Exponential families conditional conjugacy. [Bauer, 2018] **TODO** define conditional conjugacy properly

Algorithm 2: VI with conjugate family assumption

Input: A model p , variational family $q_{\phi(z)}, q_{\lambda}(z)$

```

1 while ELBO is not converged do
2   for each data point  $i$  do
3     Update  $\varphi_i \leftarrow \mathbb{E}_{\lambda} [\eta_l(\beta, x_i)]$ 
4   end
5   Update  $\lambda \leftarrow \mathbb{E}_{\varphi} [\eta_g(x, z)]$ 
6 end
```

Gradient Optimization for ELBO. We will try to solve the optimization problem from Gradient ascent perspective. This will open up opportunity for stochastic optimization [Robbins and Monro, 1951] [Robbins and Monro, 1985].

Moving from Gradient Opt to Stochastic VI

1. subsample a data point t from full data
2. use current global param λ to update local param φ_t
3. update λ

Gradient optimization step $\lambda_{t+1} = \lambda_t + \delta \nabla_{\lambda} f(\lambda_t)$. An equivalent formulation (for small $d\lambda$) is

$$\arg \max_{d\lambda} f(\lambda + d\lambda) \text{ st. } \|d\lambda\|^2 \leq \epsilon \quad (15)$$

Here we have euclidean distance metric, which is not the best choice for probability distributions. For ex - $q_{\lambda} \sim \mathcal{N}(0, 1000)$ is much closer distribution to $q_{\lambda''} \sim \mathcal{N}(10, 10000)$ than $q_{\lambda'} \sim \mathcal{N}(0, 0.001)$ is to $q_{\lambda'''} \sim \mathcal{N}(0.1, 0.001)$ even though $\|\lambda - \lambda''\| \geq \|\lambda' - \lambda'''\|$

Natural gradient of ELBO: *natural gradient* accounts for geometric structure of probability parameters (λ). They wrap the parameter space in a sensible way such that moving in same direction in different directions amounts to equal change in symmetrized KL divergence.

$$\begin{aligned} & \arg \max_{d\lambda} f(\lambda + d\lambda) \text{ st.} \\ & D_{KL}^{sym}(q_\lambda, q_{\lambda+d\lambda}) \leq \epsilon \text{ where} \\ & D_{KL}^{sym}(q, p) = KL(q||p) + KL(p||q) \end{aligned} \quad (16)$$

We need to find Riemannian metric ¹ $G(\lambda)$ which transforms euclidean distance to symmetrized KL divergence:

$$d\lambda^\top d\lambda = D_{KL}^{sym}(q_\lambda(\beta), q_{\lambda+d\lambda}(\beta)) \quad (17)$$

Using information geometry ², we can also rescale the gradients in the right space:

$$\hat{\nabla}_\lambda ELBO = G^{-1}(\lambda) \nabla_\lambda ELBO \text{ where} \quad (18)$$

$$G(\lambda) = \mathbb{E} \left[\left(\nabla_\lambda \log q_\lambda(\beta) \right) \left(\nabla_\lambda \log q_\lambda(\beta) \right)^\top \right] \quad (19)$$

$G(\lambda)$ is the Fisher information matrix. For our model class (conjugate exponential...) We've

$$\nabla_\lambda \log q_\lambda(\beta) = t(\beta) - \mathbb{E}[t(\beta)] \quad (20)$$

Combining 20 and 19

$$G(\lambda) = \nabla_\lambda^2 a(\lambda) = a''(\lambda) \quad (21)$$

From [Hoffman et al., 2013], equation of Euclidean gradient

$$\nabla_\lambda ELBO = a''(\lambda) \left(\mathbb{E}[\eta(\mathbf{x}, \mathbf{z})] - \lambda \right) \quad (22)$$

TODO refresh 22 with value in [Blei et al., 2016]

Combining 18, 22 and 21

$$\begin{aligned} g(\lambda) &= \widehat{\nabla}_\lambda ELBO = \mathbb{E}[\eta(\mathbf{x}, \mathbf{z})] - \lambda \text{ and} \\ \lambda_t &= \lambda_{t-1} + \delta_t g(\lambda_{t-1}) \\ \Rightarrow \lambda_t &= (1 - \delta_t) \lambda_{t-1} + \delta_t \mathbb{E}[\eta(\mathbf{x}, \mathbf{z})] \end{aligned} \quad (23)$$

Algorithm 3: VI with conjugate family assumption

Input: A model p , variational family $q_{\phi(z)}, q_\lambda(z)$

```

1 while ELBO is not converged do
2   for each data point i do
3     | Update  $\varphi_i \leftarrow \mathbb{E}_\lambda[\eta(\beta, x_i)]$ 
4   end
5   Update  $\lambda \leftarrow (1 - \delta_t)\lambda + \delta_t \mathbb{E}_{q(\varphi)}[\eta_g(x, z)]$ 
6 end
```

Stochastic Variational inference. in Algorithm 3 line 2-4, we have to iterate over all data to compute the new set of local variables φ . This does not scale well to large datasets. [Hoffman et al., 2013] So we have to use stochastic gradients. Noisy gradients H of f will converge to a local optimum as long as

- $\mathbb{E}[H] = \nabla f$
- Step size δ_t st: $\sum_1^\infty \delta_t = \infty$ and $\sum_1^\infty \delta_t^2 < \infty$

¹seems to be some kind of transformation

²Hope so

Now,

$$\mathbb{E}[\eta(\mathbf{x}, \mathbf{z})] = \left(\alpha_1 + \sum_1^n \mathbb{E}_q[t(z_i, x_i)], n + \alpha_2 \right)$$

Noisy gradient by sampling

1. Sample $t \sim \text{Uniform}(1, \dots, n)$
2. Rescale

$$\begin{aligned} g(\lambda) &= \left(\alpha_1 + n \mathbb{E}_q[t(z_t, x_t)], n + \alpha_2 \right) - \lambda \\ &=: \hat{\lambda} - \lambda \end{aligned}$$

Algorithm 4: Stochastic VI

Input: A model $p(\mathbf{x}, \mathbf{z})$, data \mathbf{x}
1 Initialize: variational family $q_{\phi(z)}, q_{\lambda}(z)$ with params λ_0
Result: Global variational densities $q_{\lambda}(\beta)$
2 while *Stopping criteria not met* **do**
3 Sample $t \sim \text{Uniform}(1, \dots, n)$
4 Update $\phi_t \leftarrow \mathbb{E}_{\lambda}[\eta_t(\beta, x_t)]$
5 Compute global param estimate $\hat{\lambda} = \mathbb{E}_{\varphi}[\eta_g(z_t, x_t)]$
6 Update $\lambda \leftarrow (1 - \delta_t)\lambda + \delta_t \hat{\lambda}$
7 end
8 return λ

Research on optimizing difficult variational objectives with Monte Carlo (MC) estimates. Write gradient of ELBO as expectation, compute MC estimates, use stochastic optimization with MC estimates. New approaches avoid any model-specific derivations, and are called 'Black-box' inference techniques. As examples, see - [Kingma and Welling, 2013] [Rezende et al., 2014] [Ranganath et al., 2014] [Ranganath et al., 2016] [Titsias and Lázaro-Aráoz, 2016] [Kucukelbir et al., 2017]

$$\text{ELBO} = \mathbb{E}_{q_{\nu}}[\log p_{\theta}(z, x)] - \mathbb{E}_q[\log q_{\nu}z]$$

ν params of variational family, θ params of model. We need unbiased estimates of $\nabla_{\nu, \theta} \text{ELBO}$ to maximize ELBO.

Black Box variational inference. INPROGRESS

From [Ranganath et al., 2014]

We will form the derivative of the objective as an expectation with respect to the variational approximation and then sample from the variational approximation to get noisy but unbiased gradients, which we use to update our parameters. For each sample, our noisy gradient requires evaluating the joint distribution of the observed and sampled variables, the variational distribution, and the gradient of the log of the variational distribution. This is a black box method in that the gradient of the log of the variational distribution and sampling method can be derived once for each type of variational distribution and reused for many models and applications.

We will form the ∇ELBO as an $\mathbb{E}_{q_{\lambda}}[\dots]$ and then sample S samples from the q_{λ} to get noisy but unbiased gradients (w.r.t λ), which we use to update λ . For each sample, our noisy gradient requires evaluating the $p(\mathbf{x}, \mathbf{z}_S), q(\mathbf{z}_S)$, and $\nabla \log q(\mathbf{z}_S)$. This is a black box method in that the $\nabla \log q(\mathbf{z}_S)$ and sampling method can be derived once for each type of variational distribution and reused for many models and applications.

Equation (2) of [Ranganath et al., 2014]

$$\begin{aligned} \nabla_{\lambda} \mathcal{L} &= \mathbb{E}_q \left[\nabla_{\lambda} \log q(z|\lambda) \left(\log p(x, z) - \log q(z|\lambda) \right) \right] \text{ where} \\ \mathcal{L}(\lambda) &\triangleq \mathbb{E}_{q_{\lambda}} [\log p(x, z) - \log q(z)] \text{ (ELBO)} \end{aligned} \tag{24}$$

here it says that Equation 2/3 can be derived simply using the log trick but the authors use a complicated method in paper. Also derived in [Jalil Taghia and Schn, 2018] and [Bauer, 2018]

the gradient $\nabla_\lambda \log q(z|\lambda)$ of the log of a probability distribution is called the score function or REINFORCE Basic algorithm

$$z_s \sim q(z|\lambda) \text{ for } s \in 1..S$$

$$\nabla_\lambda \mathcal{L} \approx \frac{1}{S} \sum_{s=1}^S \nabla_\lambda \log q(z_s|\lambda) \left(\log p(x, z_s) - \log q(z_s|\lambda) \right) \quad (25)$$

Rao-Blackwellization and smart Control Variates to control variance

Variance still very high. Reparameterization and amortization come to rescue (See [this](#) tutorial from David Blei)

Good notes on [Stochastic VI](#) and [Black Box VI](#) from [\[Jalil Taghia and Schn, 2018\]](#)

Reparameterization trick.

TODO

Boosting Variational inference. [\[Guo et al., 2016\]](#) **INPROGRESS**

Iterative boosting by $q_{i+1} = (1 - \gamma)q_i + \gamma h_i$. Very similar to Frank-Wolfe. Optimization is to find optimal γ and h_i at every step. γ is very similar to line search method for [\[Locatello et al., 2018\]](#) and the method is exactly same (stochastic gradient descent by taking expectations). For h_i a *Laplacian Gradient Boosting* technique is used.

Frank-Wolfe. **INPROGRESS** [\[Jaggi, 2013\]](#) [\[Pedregosa, 2018\]](#) [\[Pedregosa et al., 2018\]](#) [\[Demyanov and Rubinov, 1970\]](#)

Idea: Approximate the objective function f at iterate \mathbf{x}_t using a linear function:

$$\tilde{f}(\mathbf{s}) := f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{s} - \mathbf{x}_t \rangle$$

Find \mathbf{s} which minimizes this Linear problem (LMO) and then move in that direction by step size γ . Approximate solutions to the linear problem also suffice. Here $x, \mathcal{D} \equiv q, \mathcal{A}$

Algorithm 5: Frank-Wolfe

```

1 Constrained Optimization:  $\min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x})$ 
2  $f$  is Convex, differentiable with L-Lipschitz gradient and domain  $\mathcal{D}$  is Convex and compact
3 for  $t \in \{0, \dots, T\}$  do
4    $s^t \leftarrow \arg \min_{\mathbf{s} \in \mathcal{D}} \langle \mathbf{s}, \nabla f(\mathbf{x}^t) \rangle$ 
5    $\mathbf{x}^{t+1} \leftarrow \text{UpdateRule}(\mathbf{x}^t, s^t, t, f)$ 
6 end
```

UpdateRule can be

$$q^{t+1} \leftarrow (1 - \gamma)q^t + \gamma s^t = q^t + \gamma \overbrace{(s^t - q^t)}^{d_t} \text{ where}$$

$$\text{Variant0 : } \gamma \leftarrow \frac{2}{t+2} \quad (26)$$

$$\text{Variant1 : } \gamma \leftarrow \arg \min_{\gamma \in [0,1]} f((1 - \gamma)q^t + \gamma s^t) \quad (27)$$

$$g_t \leftarrow -\langle \nabla f(\mathbf{x}_t), d_t \rangle$$

Exitcondition : $g_t < \delta$

$$\text{Variant2 : } \gamma \leftarrow \min \left(\frac{g_t}{L \|d_t\|^2}, 1 \right) \quad (28)$$

Variant3 :

$$q^{t+1} \leftarrow \arg \min_{q \in \bigcup_{i=1}^t s^i} f(q) \quad (29)$$

$$(30)$$

Constraint \mathcal{D}	LMO problem
norm $\ x\ \leq 1$	$-\partial \ \cdot\ _*$ Subgradients of corresponding dual norm
l_1 norm $\ x\ _1 \leq 1$	$-\partial \ \nabla f(\mathbf{x}_t)\ _\infty$
Trace norm $\ X\ _{tr} \leq 1$	Operator norm $s_t \in -\ \nabla f(X_t)\ _{op}$
sum of singular values	Largest singular value

Table 1: LMO problem for well known constraints

$$\gamma \leftarrow \min \left\{ \frac{g_t}{L \text{diam}(\mathcal{D})^2}, 1 \right\}$$

Frank-Wolfe Convergence. INPROGRESS

$$\begin{aligned} \overbrace{g_t}^{\text{Gap}} &:= \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{s}_t \rangle \\ &\geq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle \\ &\geq \underbrace{f(\mathbf{x}_t) - f(\mathbf{x}^*)}_{\epsilon_t} \quad \text{Convexity} \end{aligned} \tag{31}$$

$$\tag{32}$$

Using quadratic upper bound 3 L-continous gradient can be relaxed to this, We get,

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2 \tag{33}$$

$$\begin{aligned} \Rightarrow f(\mathbf{x}_{t+1}) &= f\left((1-\gamma)\mathbf{x}_t + \gamma\mathbf{s}_t\right) \leq f(\mathbf{x}_t) - \gamma g_t + \frac{L\gamma^2}{2} \|\mathbf{s}_t - \mathbf{x}_t\|^2 \\ &\leq f(\mathbf{x}_t) - \gamma\epsilon_t + \frac{L\gamma^2}{2} \|\mathbf{s}_t - \mathbf{x}_t\|^2 \\ \Rightarrow \epsilon_{t+1} &\leq (1-\gamma)\epsilon_t + \frac{L\gamma^2}{2} \mathcal{D}_t^2 \end{aligned} \tag{34}$$

$$\leq (1-\gamma_t)\epsilon_t + \gamma_t^2 C \left(= \frac{L}{2} \text{diam}(\mathcal{D}) \right) \tag{35}$$

$$\tag{36}$$

Goal: To show, for $\gamma_t = \frac{2}{t+2}$

$$\epsilon_t \leq \frac{4C}{t+2} \tag{37}$$

Using induction, for $t = 0$, $\epsilon_0 = C \leq \frac{4C}{2}$. At step t

$$\begin{aligned} \epsilon_{t+1} &\leq \left(1 - \frac{2}{t+2}\right)\epsilon_t + \left(\frac{2}{t+2}\right)^2 C \\ &\leq \left(\frac{t}{t+2}\right) \cdot \left(\frac{4C}{t+2}\right) + \left(\frac{1}{t+2}\right) \cdot \left(\frac{4C}{t+2}\right) \\ &= \frac{4C}{t+2} \frac{t+1}{t+2} \\ &\leq \frac{4C}{t+2} \frac{t+2}{t+3} \\ &= \frac{4C}{t+1+2} \end{aligned}$$

Revisiting Frank-Wolfe: Projection Free Sparse Convex Optimization. [Jaggi, 2013] **INPROGRESS** Variants of FW in 6

- **Approximate LMO:** $\epsilon := \frac{1}{2}\delta\gamma C_f$ additive approximate error
- **Fully Corrective:** in Equation 29, is only a degenerate boosting as only one atom \mathbf{s} is chosen at each time. If we change the search space to

$$q^{t+1} \leftarrow \arg \min_{q \in \text{conv}(\bigcup_{i=1}^t \mathbf{s}^i)} f(q)$$

Then the progress made per iteration would be more but the search problem would not be much easier than the original problem

Curvature Constant C_f of a convex and differentiable f :

$$C_f := \sum_{\mathbf{x}, \mathbf{s} \in \mathcal{D}, \gamma \in [0,1], \mathbf{y} = \mathbf{x} + \gamma(\mathbf{s} - \mathbf{x})} \frac{2}{\gamma^2} \left(f(\mathbf{y}) - f(\mathbf{x}) - \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle \right)$$

Note that $C_f = \frac{2}{\gamma^2} (f - \tilde{f})$ means that for bounded C_f , deviation of f from \tilde{f} will also be bounded. $f - \tilde{f}$ is also called *Bregman divergence*. If ∇f is L-Lipschitz continuous on \mathcal{D} w.r.t some norm $\|\cdot\|$, then $C_f \leq \text{diam}_{\|\cdot\|}(\mathcal{D}^2)L$

Boosting Black Box Variational inference. **INPROGRESS**

Boosting introduced in [Guo et al., 2016], connection with FW in [Locatello et al., 2017]. define a Linear Minimization Problem (LMO) as $\mathbf{LMO}_{\mathcal{A}}(y) := \arg \min_{s \in \mathcal{A}} \langle y, s \rangle$ In line 3 of 5, rewrite it as

$$s^t \leftarrow (\delta - \text{Approx-})\mathbf{LMO}_{\mathcal{A}}(\nabla f(q^t))$$

Algorithm for LMO in section 4 of [Locatello et al., 2018]. In Theorem 2, Curvature $\mathcal{C}_{f,\mathcal{A}}$ is bounded for D^{KL} if param. space of densities in \mathcal{A} is bounded. In section 3, a bounded curvature for D^{KL} is obtained.

Black box LMO:

In this case $f(q^t) = \text{KL}(q^t(\mathbf{z})||p(\mathbf{x}, \mathbf{z}))$. Assuming θ are the parameters defining variational family $\mathcal{Q} \equiv \mathcal{A}$ We've to find $\nabla_{\theta} f(q^t)$, more specifically, we've to find

$$s^t \leftarrow (\delta - \text{Approx.}) \arg \min_{s \in \mathcal{A}} \langle \nabla \text{KL}(q^t(\mathbf{z})||p(\mathbf{x}, \mathbf{z})), s \rangle$$

TODO add how [Guo et al., 2016] [Locatello et al., 2017] deal with optimization of LMO. Also add the part about $\text{conv}(\mathcal{A})$ being sufficient instead of \mathcal{A} .

Convergence of SGD not fully understood. To guarantee convergence of FW, solution of LMO should not be degenerate. This translates to a constraint on $\|s\|_{\infty}$ which is not practical. Every pdf with bounded $\|\cdot\|_{\infty}$ has bounded entropy and the converse holds true in most cases of interest. (Gaussian, Laplacian, ...). Assume \mathcal{A} is such a family and $\bar{\mathcal{A}}$ is \mathcal{A} w/o l_{∞} norm constraint. **TODO** ask

$$\arg \min_{s \in \bar{\mathcal{A}}, \mathcal{H}(s) \geq -M} \langle \nabla \text{KL}(q^t(\mathbf{z})||p(\mathbf{x}, \mathbf{z})), s \rangle \stackrel{?}{=} \arg \min_{s \in \bar{\mathcal{A}}, \mathcal{H}(s) \geq -M} \left\langle s, \log \frac{q^t}{p} \right\rangle$$

Using Lagrange multiplier λ

$$\begin{aligned} & \left\langle s, \log \left(\frac{s}{\sqrt{\frac{p}{q^t}}} \right) \right\rangle \\ & \equiv \arg \min_{s \in \bar{\mathcal{A}}} \text{KL} \left(s || \sqrt{\frac{p}{q^t}} Z \right) \\ \text{RELBO}(s, \lambda) &:= \mathbb{E}_s [\log p] - \mathbb{E}_s [\log q^t] - \lambda \mathbb{E}_s [\log s] \end{aligned} \tag{38}$$

For true LMO solution, will need to maximize for λ . Might end in saddle, fix or slowly decrease with time $\frac{1}{\sqrt{t+1}}$

6 Ideas

6.1 Line search

Line search in [27](#) is not working very well.

line search

$$\gamma' = \argmin_{\gamma \in [0,1]} KL \left(\underbrace{q^t + \gamma(s - q^t)}_{q_r^t} \parallel p \right)$$

$$s = \argmin KL \left(s \parallel \sqrt{\frac{p}{q_r}} z \right)$$

RELBO
CLR

grad.

$$\nabla_{\gamma} = \int q_r^t \log \left(\frac{q_r^t}{p} \right) dz$$

$$= - \int \underbrace{\nabla_{\gamma} q_r^t}_{s - q^t} \log \left(\frac{q_r^t}{p} \right) + \underbrace{q_r^t}_{\frac{1}{q_r^t}} \cdot \underbrace{\nabla_{\gamma} \log \left(\frac{q_r^t}{p} \right)}_{\frac{1}{q_r^t} \cdot \nabla_{\gamma} q_r^t}$$

$$= - \int \nabla_{\gamma} q_r^t \left(1 + \log q_r^t - \log p \right)$$

$$= - \int (s - q^t) \left(1 + \log q_r^t - \log p \right)$$

$$= - \mathbb{E}_s \left[1 + \log q_r^t - \log p \right] + \mathbb{E}_{q^t} \left[1 + \log q_r^t - \log p \right]$$

res-s ← res-q ←

$$= \begin{cases} \mu^{t+1} \leftarrow [\mu^t, w_s] \\ \sigma^{t+1} \leftarrow [\sigma^t, \delta_s] \end{cases}$$

$$\gamma \leftarrow \frac{2}{k+2}$$

$$w = ((1-\gamma)w, \gamma)$$

q_{next}

$$new\ mix = (1-\gamma)q + \gamma s$$

$$\hat{\nabla}_{\gamma} \approx \frac{1}{k} \frac{1}{\text{var } N} (res - q - res - s)$$

$$\gamma_{new} = \gamma_{old} + \alpha \cdot \frac{1}{t+1} (\hat{\nabla}_{\gamma})$$

once 0
will always
report to 0!!

project onto [0,1]

```

1 def line_search_dkl(weights, locs, diags, mu_s, cov_s, x, k):
2     """Perform line search for the best step size gamma.
3
4     Uses gradient ascent to find gamma that minimizes
5     KL(q_t + gamma (s - q_t) || p)
6
7     Args:
8         weights: weights of mixture components of q_t
9         locs: means of mixture components of q_t
10        diags: deviations of mixture components of q_t
11        mu_s: mean for LMO Solution s
12        cov_s: cov matrix for LMO solution s
13        x: target distribution p
14        k: iteration number of Frank-Wolfe
15    Returns:
16        Computed gamma
17    """
18    def softmax(v):
19        return np.log(1 + np.exp(v))
20    # no. of samples to approximate  $\nabla_\gamma$ 
21    N_samples = 10
22    # Create current iter q_t
23    weights = [weights]
24    qt_comps = [
25        Normal(
26            loc=tf.convert_to_tensor(locs[i]),
27            scale=tf.convert_to_tensor(diags[i])) for i in range(len(locs))
28    ]
29    qt = Mixture(
30        cat=Categorical(probs=tf.convert_to_tensor(weights)),
31        components=qt_comps,
32        sample_shape=N)
33    qt = InfiniteMixtureScipy(stats.multivariate_normal)
34    qt.weights = weights[0]
35    qt.params = list(
36        zip([[l] for l in locs], [[softmax(np.dot(d, d))] for d in diags]))
37    # samples from q_t
38    sample_q = qt.sample_n(N_samples)
39    # create and sample from s
40    s = stats.multivariate_normal([mu_s],
41                                  np.dot(np.array([cov_s]), np.array([cov_s])))
42    sample_s = s.rvs(N_samples)
43    # q_{t+1} is mixture of q_t and s with weights (1 - \gamma) and \gamma
44    # Set its corresponding parameters and weights
45    new_locs = copy.copy(locs)
46    new_diags = copy.copy(diags)
47    new_locs.append([mu_s])
48    new_diags.append([cov_s])
49    # initialize \gamma
50    gamma = 2. / (k + 2.)
51    # no. steps of gradient ascent
52    n_steps = 10
53    prog_bar = ed.util.Progbar(n_steps)
54    for it in range(n_steps):
55        print("line_search iter %d, %.5f" % (it, gamma))
56        new_weights = copy.copy(weights)
57        new_weights[0] = [(1. - gamma) * w for w in new_weights[0]]
58        new_weights[0].append(gamma)
59        # create q_{t+1}^gamma
60        q_next = InfiniteMixtureScipy(stats.multivariate_normal)
61        q_next.weights = new_weights[0]
62        q_next.params = list(
63            zip([[l] for l in new_locs], [[np.dot(d, d)] for d in new_diags]))
64        # Computes  $\mathbb{E}[\dots] \propto \sum_v \log p - \log q_{t+1}^\gamma$ 
65        def px_qx_ratio_log_prob(v):
66            Lambda = 1.

```

```

67     ret = x.log_prob([v]).eval()[0] - q_next.log_prob(v)
68     ret /= Lambda
69     return ret
70     # Samples w.r.t s
71     rez_s = [
72     px_qx_ratio_log_prob(sample_s[ss]) for ss in range(len(sample_s))
73     ]
74     # Samples w.r.t qt+1
75     rez_q = [
76     px_qx_ratio_log_prob(sample_q[ss]) for ss in range(len(sample_q))
77     ]
78     # Gradient ascent step, step size decreasing as  $\frac{1}{it+1}$ 
79     gamma = gamma + 0.1 * (sum(rez_s) - sum(rez_q)) / (N_samples *
80     (it + 1.))
81     # Projecting it back to [0, 1], too small range?
82     # FIXME(sauravshekhar) if projected to 0, all iterations will be same?
83     if gamma >= 1 or gamma <= 0:
84     gamma = max(min(gamma, 1.), 0.)
85     break
86     return gamma

```

changes for measuring variance of $\mathbb{E}_s[\cdot]$ and $\mathbb{E}_{q_{t+1}}[\cdot]$.

```

1     ...
2     grad_gamma = []
3     for it in range(n_steps):
4     ...
5     # Samples w.r.t s
6     rez_s = np.asarray([
7     px_qx_ratio_log_prob(sample_s[ss]) for ss in range(len(sample_s))
8     ])
9     # Samples w.r.t qt+1
10    rez_q = np.asarray([
11    px_qx_ratio_log_prob(sample_q[ss]) for ss in range(len(sample_q))
12    ])
13    grad_gamma.append({'E_s': rez_s, 'E_q': rez_q, 'gamma': gamma})
14    ...
15    # Write grad_gamma to outdir/line_search_samples_<n_samples>.npy.<fw_iter>

```

Metrics on original version.

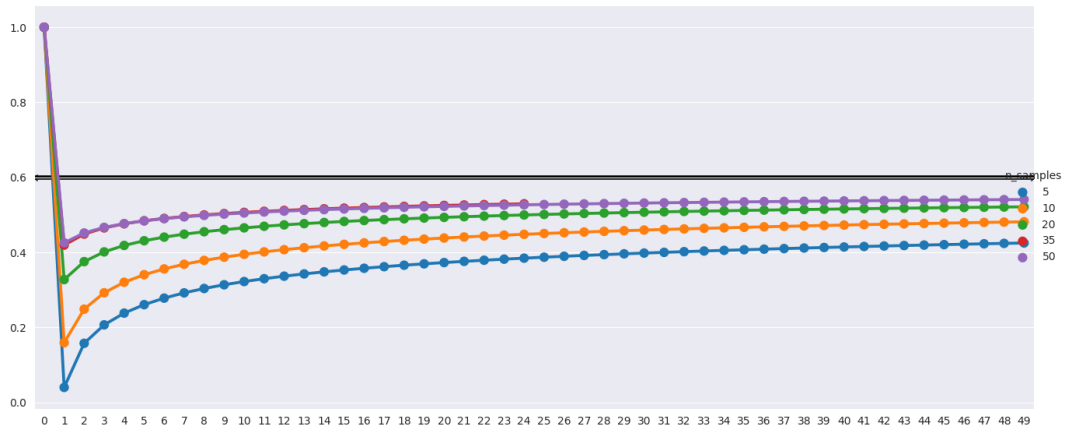


Figure 1: gamma with iterations for different n_samples

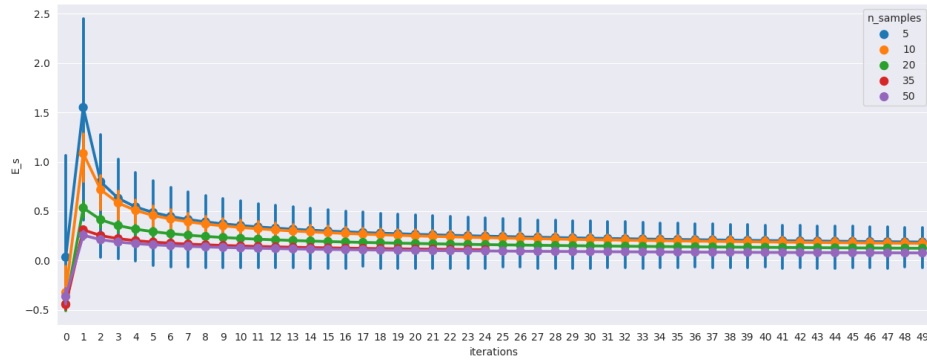


Figure 2: E_s with different $n_samples$

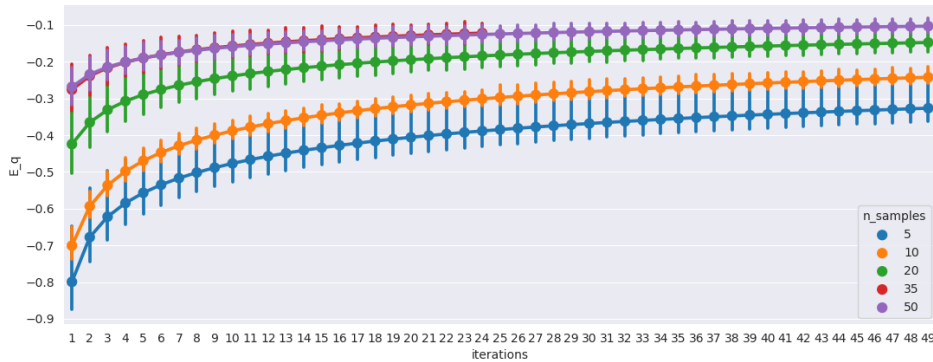


Figure 3: E_q with different $n_samples$

?? begins with iter 1 as iter 0 has very high variance

6.2 Measuring smoothness

INPROGRESS Computing optimal γ directly from eqn 1 of [Pedregosa et al., 2018]

$$\begin{aligned}
 f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) + \gamma \langle \nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{x}_t \rangle + \frac{\gamma^2}{2} L_t \|\mathbf{s}_t - \mathbf{x}_t\|^2 \\
 \Rightarrow L_t &\geq \frac{f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) + \gamma \langle \nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{x}_t \rangle}{\frac{\gamma^2}{2} \|\mathbf{s}_t - \mathbf{x}_t\|^2} \\
 &= \frac{f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) + \gamma \langle \nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{x}_t \rangle}{\frac{\gamma^2}{2} \text{KL}(\mathbf{s}_t || \mathbf{q}_t)} \\
 \langle f, g \rangle &= \int f(\theta) f(\theta) d\theta \quad // \text{functional dot product}
 \end{aligned}$$

Code changes.

```

1  def kl(mu_q, std_q, mu_p, std_p):
2      ...
3
4  def func_dot_product():
5      ...

```

6.3 Other optimization algorithm

TODO As shown in link 3, Frank-Wolfe converges slower than Projected Gradient Descent in Practice. See [Locatello et al., 2017] to see why we use FW and if it can be replaced. (will have to derive new convergence proofs and boosting won't be as integrated into the optimization algorithm as before).

References

- [Balabdaoui and van de Geer, 2016] Balabdaoui, F. and van de Geer, S. (2016). Fundamentals of mathematical statistics.
- [Bauer, 2018] Bauer, S. (2018). Probabilistic graphical models for image analysis.
- [Blei et al., 2016] Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. (2016). Variational inference: A review for statisticians. *arXiv preprint arXiv:1601.00670*.
- [Demjanov and Rubinov, 1970] Demjanov, V. F. and Rubinov, A. M. (1970). Approximate methods in optimization problems. (*Modern Analytic and Computational Methods in Science and Mathematics. IX*).
- [Guo et al., 2016] Guo, F., Wang, X., Fan, K., Broderick, T., and Dunson, D. B. (2016). Boosting variational inference. *arXiv preprint arXiv:1611.05559*.
- [Hoffman et al., 2013] Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. (2013). Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347.
- [Jaggi, 2013] Jaggi, M. (2013). Revisiting frank-wolfe: Projection-free sparse convex optimization. In *ICML (1)*, pages 427–435.
- [Jalil Taghia and Schn, 2018] Jalil Taghia, L. M. and Schn, T. (2018). Probabilistic machine learning.
- [Kingma and Welling, 2013] Kingma, D. P. and Welling, M. (2013). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- [Kucukelbir et al., 2017] Kucukelbir, A., Tran, D., Ranganath, R., Gelman, A., and Blei, D. M. (2017). Automatic differentiation variational inference. *The Journal of Machine Learning Research*, 18(1):430–474.
- [Locatello et al., 2018] Locatello, F., Dresdner, G., Khanna, R., Valera, I., and Rätsch, G. (2018). Boosting black box variational inference. *arXiv preprint arXiv:1806.02185*.
- [Locatello et al., 2017] Locatello, F., Khanna, R., Ghosh, J., and Rätsch, G. (2017). Boosting variational inference: an optimization perspective. *arXiv preprint arXiv:1708.01733*.
- [Pedregosa, 2018] Pedregosa, F. (2018). Notes on the frank-wolfe algorithm, part i. [Online; posted 21-March-2018].
- [Pedregosa et al., 2018] Pedregosa, F., Askari, A., Negiar, G., and Jaggi, M. (2018). Step-size adaptivity in projection-free optimization. *arXiv preprint arXiv:1806.05123*.
- [Ranganath et al., 2014] Ranganath, R., Gerrish, S., and Blei, D. (2014). Black box variational inference. In *Artificial Intelligence and Statistics*, pages 814–822.
- [Ranganath et al., 2016] Ranganath, R., Tran, D., and Blei, D. (2016). Hierarchical variational models. In *International Conference on Machine Learning*, pages 324–333.
- [Rezende et al., 2014] Rezende, D. J., Mohamed, S., and Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, pages 1278–1286.
- [Robbins and Monro, 1951] Robbins, H. and Monro, S. (1951). A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3):400–407.
- [Robbins and Monro, 1985] Robbins, H. and Monro, S. (1985). A stochastic approximation method. In *Herbert Robbins Selected Papers*, pages 102–109. Springer.
- [Titsias and Lázaro-Gredilla, 2014] Titsias, M. and Lázaro-Gredilla, M. (2014). Doubly stochastic variational bayes for non-conjugate inference. In *International Conference on Machine Learning*, pages 1971–1979.