

1 Ideas

1.1 Line search

Line search in ?? is not working very well.

line search

$$r' = \underset{r \in [0,1]}{\operatorname{argmin}} \quad \underset{\tilde{D}_{KL}(r)^\uparrow}{KL(q^t + r(s-q^t) \| p)} = \nabla_r f(\tilde{q}_t)$$

grad descent,

$$r \leftarrow r - \eta \nabla_r \tilde{D}_{KL}(r)$$

func. grad,

$$\nabla_{\tilde{q}} f(\tilde{q})$$

$$\nabla_r KL(q^t + r(s-q^t) \| p)$$

$$= \log \frac{r}{p}$$

$$\int \underbrace{(q^t + r(s-q^t))}_{\tilde{q}_t} \log \left(\frac{q^t + r(s-q^t)}{p} \right) dr$$

$$-g_n = \langle \nabla f(n), s \rangle - \langle \nabla f(n), n \rangle$$

$$= \nabla_r \int \tilde{q}_t \log \frac{\tilde{q}_t}{p}$$

$$= \int \left(\nabla_r \tilde{q}_t \right) \log \frac{\tilde{q}_t}{p} + \cancel{\tilde{q}_t \nabla_r \log \frac{\tilde{q}_t}{p}}$$

\downarrow
 \perp
 $s - q^t$

\downarrow
 $\frac{1}{\tilde{q}_t} \nabla_r \tilde{q}_t$
 \downarrow
 $s - q^t$

$$= \int (s - q^t) \log \frac{\tilde{q}_t}{p} + (s - q^t)$$

$$= \int (s - q^t) \left(1 + \log \frac{\tilde{q}_t}{p} \right)$$

$$= \int s \log \frac{\tilde{q}_t}{p} - \int q^t \log \frac{\tilde{q}_t}{p}$$

$$= \langle s, \log \frac{\tilde{q}_t}{p} \rangle - \langle q^t, \log \frac{\tilde{q}_t}{p} \rangle$$

$$= \mathbb{E}_s [\log \tilde{q}_t - \log p] - \mathbb{E}_{q^t} [\log \tilde{q}_t - \log p]$$

$$= \mathbb{E}_s [\nabla_{\tilde{q}_t} f(\tilde{q}_t)] - \mathbb{E}_{q^t} [\nabla_{q^t} f(\tilde{q}_t)]$$

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 q_{t+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.

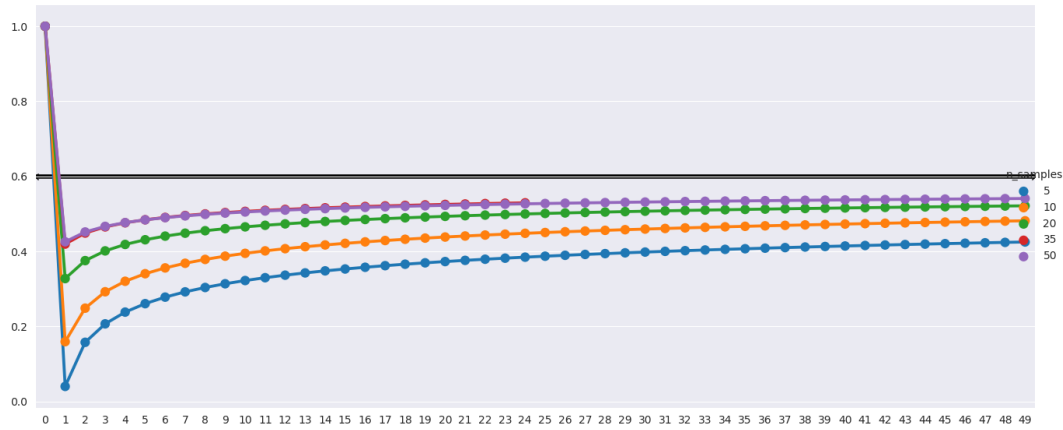


Figure 1: gamma with iterations for different n_samples

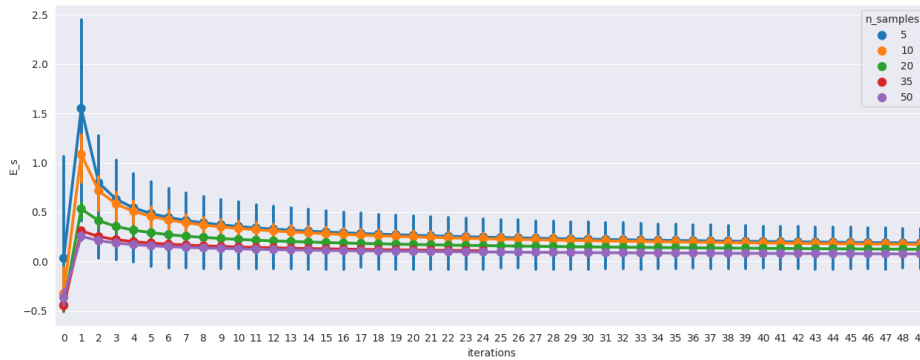


Figure 2: E_s with different n_samples

values have changed slightly after moving everything to Tensorflow/Edward Remake plots?

1.2 Adaptive Frank-Wolfe from smoothness estimators

Algorithm 1 of [Pedregosa et al., 2018]

Code changes.

1.3 Measuring smoothness

in progress

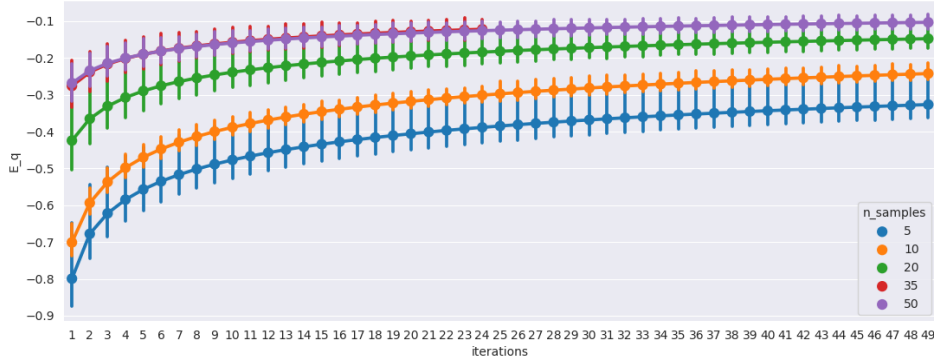


Figure 3: E_q with different $n_samples$
 ?? begins with iter 1 as iter 0 has very high variance

Algorithm 1: Adaptive Frank-Wolfe for Boosting BBVI

Input: $q_0 \in \mathcal{D}$, initial Lipschitz estimate L_{-1} , line search parameters $\tau > 1, \eta \in (0, 1]$

```

1 for  $t = 1 \dots T$  do
2    $s_t \leftarrow LMO_{\mathcal{A}}(\nabla f(q^t))$ 
3    $g_t \leftarrow \langle \nabla f(q_t), q_t \rangle - \langle \nabla f(q_t), s_t \rangle$  // Gap  $\geq 0$ 
4   Find smallest integer  $i$  s.t
5    $f(q_t + \gamma_t(s_t - q_t)) \leq Q_t(\gamma_t)$  Where
6    $Q_t(\gamma) := f(q_t) - \gamma g_t + \frac{\gamma^2 L_t}{2} d(s_t, q_t)$  // Quadratic upper bound ??
7    $L_t \leftarrow \tau^i \eta L_{t-1}$  and  $\gamma_t \leftarrow \min\left(\frac{g_t}{L_t d(s_t, q_t)}, 1\right)$ 
8 end
9 return  $q_T$ 
```

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

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

Code changes.

```

1 def grad_kl(q, p, theta):
2     # Functional Gradient w.r.t q  $\nabla KL(q||p) = \log q - \log p$ 
3     ...
4
5
6 def lmo(y, p, ...):
7     #  $f = \mathcal{D}^{KL}(y||p)$ 
8     #  $\langle \nabla f, y \rangle = \mathbb{E}_y \nabla f$ 
9     ...
```

1.4 Other optimization algorithm

todo

As shown in link ??, 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).

1.5 Entropy Regularization and Noise addition using Optimal Transport

In LMO, [Locatello et al., 2018] uses Entropy Regularization in place of norm constrained optimization. It can be replaced with something simpler See [Tolstikhin et al., 2017] [Dong Liu, 2018] [Bernton et al., 2017] [Jordan et al., 1998] [here](#) And [Peyré et al., 2017] part 4.

1.6 Sensibility of Distribution

Given observable $r(x)$ with mean μ , Probability distribution $p(x)$ that satisfies the constraint $\mathbb{E}_{p(x)}[r(x)] = \mu$ is a function of parameter μ , $p(x) \equiv p(x, \mu)$. Sensibility of a distribution w.r.t μ is defined as the expected squared deviation given by the perturbation of μ . See slides of lecture 2 of SLT

$$\mathcal{S}(p) = \mathbb{E} \left[\left(\frac{\partial_\mu p}{p} \right)^2 \right]$$

1.7 Port code to Tensorflow Probability

see ?? . Issue is if `tfp.edward2` will have support for Variational Inference and ELBO etc.

References

- [Bernton et al., 2017] Bernton, E., Jacob, P. E., Gerber, M., and Robert, C. P. (2017). Inference in generative models using the wasserstein distance. *arXiv preprint arXiv:1701.05146*.
- [Dong Liu, 2018] Dong Liu, Minh Thnh Vu, S. C. L. K. R. (2018). Entropy-regularized optimal transport generative models. *arXiv preprint arXiv:1811.06763*.
- [Jordan et al., 1998] Jordan, R., Kinderlehrer, D., and Otto, F. (1998). The variational formulation of the fokker–planck equation. *SIAM journal on mathematical analysis*, 29(1):1–17.
- [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 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*.
- [Peyré et al., 2017] Peyré, G., Cuturi, M., et al. (2017). Computational optimal transport. Technical report, École Normale Supérieure.
- [Tolstikhin et al., 2017] Tolstikhin, I., Bousquet, O., Gelly, S., and Schoelkopf, B. (2017). Wasserstein auto-encoders. *arXiv preprint arXiv:1711.01558*.