# 1 Ideas

## 1.1 Line search

Line search in **??** is not working very well.

line Search

$$\gamma' = \underset{\gamma \in [0,1]}{\text{argmin}} \quad KL\left(q^t + \gamma(s - q^t) \| P\right) \qquad = \nabla_\gamma f(\tilde{q}_t)$$

$$\underbrace{\tilde{D}_{KL}(\gamma)}$$

grad descent,

$$\gamma \leftarrow \gamma - \eta \, \nabla_\gamma \tilde{D}_{KL}(\gamma)$$

func. grad,

$$\nabla_{n_q} f(n)$$

$$= \ln \frac{n}{P}$$

$$\nabla_\gamma KL\left(q^t + \gamma(s - q^t) \| P\right)$$

$$\downarrow$$

$$\int \underbrace{(q^t + \gamma(s - q^t))}_{\tilde{q}_t} \ln\left(\frac{q^t + \gamma(s - q^t)}{P}\right) d\alpha \qquad -g_n = \langle \nabla f(n), s \rangle$$
$$-\langle \nabla f(\alpha), n \rangle$$

$$= \nabla_\gamma \int \tilde{q}_t \ln \frac{\tilde{q}_t}{P}$$

$$= \int \left(\nabla_\gamma \tilde{q}_t\right) \ln \frac{\tilde{q}_t}{P} \; + \; \tilde{q}_t \, \nabla_\gamma \ln \frac{\tilde{q}_t}{P}$$

$$\underset{s - q^t}{\Big\downarrow} \qquad\qquad \frac{1}{\tilde{q}_t} \nabla_\gamma \hat{q}_t$$

$$\searrow (s - q^t)$$

$$= \int (s - q^t) \ln \frac{\tilde{q}_t}{P} \; + \; (s - q^t)$$

$$= \int (s - q^t) \left(1 + \ln \frac{\check{q}_t}{P}\right)$$

$$= \int s \ln \frac{\tilde{q}_t}{P} \; - \; \int q^t \ln \frac{\tilde{q}_t}{P}$$

$$= \langle s, \ln \frac{\tilde{q}_t}{P} \rangle \; - \; \langle q^t, \ln \frac{\check{q}_t}{P} \rangle$$

$$= \mathbb{E}_s\left[\ln \tilde{q}_t - \ln P\right] \; - \; \mathbb{E}_{q^t}\left[\ln \check{q}_t - \ln P\right]$$

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

```python
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 ∇_γ
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 − γ) and γ
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 γ
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}^γ
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 E[...] ∝ Σ_v log p − log q_{t+1}^γ
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 q_{t+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 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 $\mathbb{E}_s\left[\cdot\right]$ and $\mathbb{E}_{q_{t+1}^\gamma}\left[\cdot\right]$.

```
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>
```
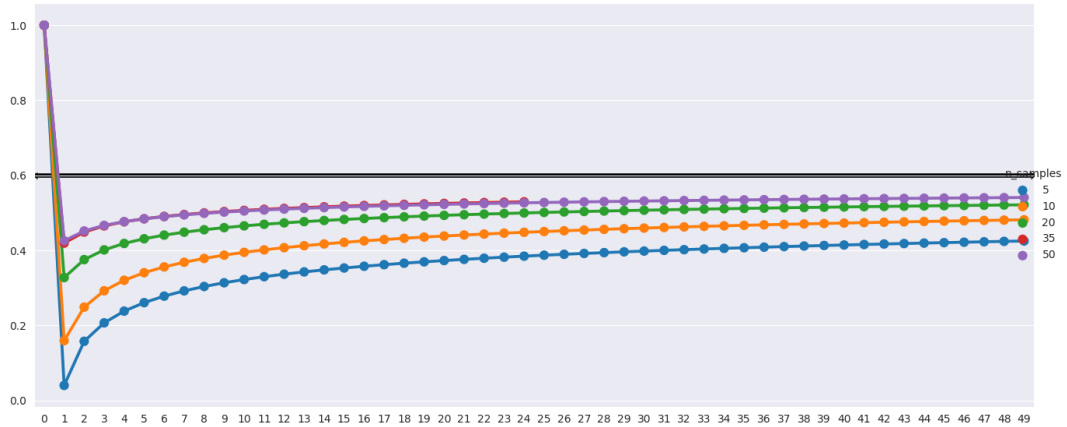
Metrics on original version.



Figure 1: gamma with iterations for different n_samples

## 1.2   Adaptive Frank-Wolfe from smoothness estimators

Algorithm 1 of [Pedregosa et al., 2018]
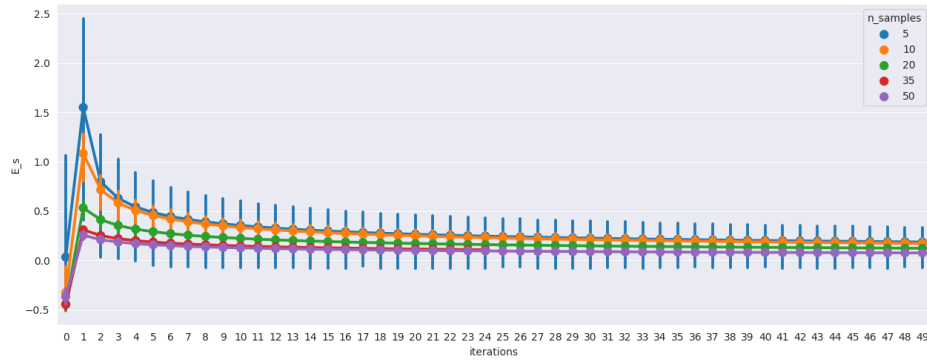
**Code changes.**
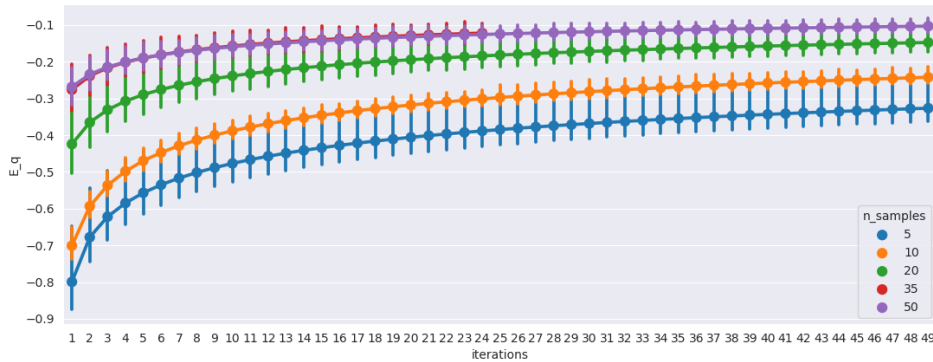
4

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

## 1.3 Measuring smoothness

in progress

Computing optimal $\gamma$ directly from eqn 1 of [Pedregosa et al., 2018]

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \gamma\langle\nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{x}_t\rangle + \frac{\gamma^2}{2}L_t||\mathbf{s}_t - \mathbf{x}_t||^2$$

$$\Rightarrow L_t \geq \frac{f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) + \gamma\langle\nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{x}_t\rangle}{\frac{\gamma^2}{2}||\mathbf{s}_t - \mathbf{x}_t||^2}$$

$$\frac{f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) + \gamma\langle\nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{x}_t\rangle}{\frac{\gamma^2}{2}\mathrm{KL}\left(s_t||q_t\right)}$$

**Code changes.**

```python
def kl(mu_q, std_q, mu_p, std_p):
    ...

def func_dot_product():
    ...
```

## 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).

5

**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** $t = 1 \ldots T$ **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** $q_T$

## 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 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*.