$\pi \propto \exp(-V)$

MCMC: $X_{(1,\dots,)} X_T \sim \hat{\pi} \approx \pi$

$\hat{\mu}_\pi = \frac{1}{T} \sum X_t$, $\hat{\Sigma}_\pi = \frac{1}{T} \sum (X_t - \hat{\mu}_\pi) \otimes (X_t - \hat{\mu}_\pi)$

MCMC as discrete Langevin

$$dX_t = -\nabla V(X_t) dt + \sqrt{2} dB_t$$

has $\pi$ as unique stationary dist.

$V$ is strongly convex $\Rightarrow$
appoximately sample $\pi$ after
$O(d)$ queries to $\nabla V$

---

Variational inference:

$\hat{\pi} = \arg \min_{p \in P} KL(p \| \pi)$

What is $P$

(1) $P_2(\mathbb{R}^d)$ — measures with bounded second moments

(2) $P$ is a product measure - mean-field
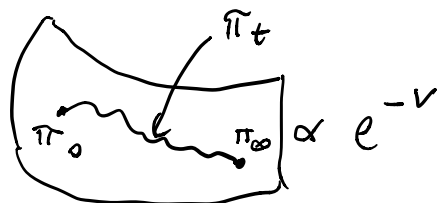
(3) $P$ is a Gaussian

How does $\hat{\pi}$ approximate $\pi$?

$$(m, \Sigma) \longmapsto KL\left(\mathcal{O}_{(m, \Sigma)} \| \pi\right) \quad \text{not convex}$$

$$(m, \Sigma) \longmapsto KL\left(\mathcal{O}_{(m, \Sigma)} \| \pi\right) \quad \begin{array}{l} \text{not convex} \\ \text{in Fisher} \\ \text{geometry} \end{array}$$

---

Wasserstein gradient flows.

$(\pi_t)_{t \geq 0}$ is the marginal dist. of Langevin diffusion.

$$\pi_t = \text{Law}(X_t) \quad , \quad dX_t = -\nabla V(X_t) dt + \sqrt{2} \, dB_t$$

$$\partial_t \pi_t = -\text{div}(\pi_t \nabla V) + \Delta \pi_t \qquad \begin{array}{l} \text{Fokker-} \\ \text{Planck} \end{array}$$



$\pi_\infty \propto e^{-V}$

$$W_2(u,v) = \left[ \inf_{\gamma \in C(u,v)} \|x-y\|^2 \, d\gamma(x,y) \right]^{1/2}$$

$(\mathcal{P}_2(\mathbb{R}^d), W_2)$ is a metric space.

thm (JKO '98):

The law $(\pi_t)_{t\geq 0}$ of the Langevin diffusion is a gradient flow of $KL(\cdot \| \pi)$ on the Wasserstein space $(\mathcal{P}_2(\mathbb{R}^d), W_2)$.

$$\partial_t(\pi_t) = -\text{div}(\pi_t \nabla v) + \Delta \pi_t \iff \dot{X}_t = -\nabla_{W_2} KL(\pi_t \| \pi)(X_t)$$

$$\dot{X}_t = -\nabla v(X_t) - \nabla \log \pi_t(X_t)$$

JKO-scheme
$$\pi_{n,k+1} := \arg\min_{u \in \mathcal{P}_2(\mathbb{R}^d)} \left[ KL(u \| \pi) + \frac{1}{2h} W_2^2(u, u_{n,k}) \right]$$

how: $\pi_t$ is unknown so how do you implement this?

Idea: evolve N particles $X_t^{(1)}, \dots, X_t^{(N)}$ and some how take an expection.

$$K(x,y) = \begin{bmatrix} e^{-\|x_1 - y_1\|^2/\sigma^2} \\ \vdots \\ e^{-\|x_d - y_d\|^2/\sigma^2} \end{bmatrix}$$

$$\dot{X}_t = \nabla_{W_2} KL(\pi_t \| \pi)(X_t)$$
$$= -\nabla V(X_t) - \nabla \log \pi_t(X_t)$$

$$\nabla \log \pi_t(x) \approx \int K(x,y) \nabla \log \pi_t(y) d\pi_t(y)$$

$$= \int K(x,y) \frac{\nabla \pi_t(y)}{\pi_t(y)} \pi_t(y) dy$$

$$= \int k(x,y) \nabla \pi_t(y) dy$$

$$= -\int \nabla_y k(x,y) \pi_t(y) dy$$
$$= -\mathbb{E}_{x_t \sim \pi_t}[\nabla_y K(x, X_t)]$$

Project Wasserstein gradient onto RKHS.

$$\nabla \log \pi_t(x) \approx -\frac{1}{N} \sum_{i=1}^{N} \nabla_y K(x, x_t^{(i)})$$

$$\dot{X}_t^{(j)} = -\nabla V(X_t^{(i)}) + \frac{1}{N} \sum_{i=1}^{N} \nabla_y K(X_t^{(j)}, X_t^{(i)})$$

$$, j = 1, \ldots, N$$

Some Issues:
1) K?
2) Does not scale for $d > 3$
3) VI with unclear $p$

Will come back to this

$(\pi_t)_{t \geq 0}$

$$\pi_t = Law(X_t) \qquad dX_t = -\nabla V(X_t)dt + \sqrt{2} \, dB_t$$

$$m_t = \mathbb{E}(X_t), \qquad \Sigma_t = cov(X_t)$$

$$\dot{m}_t = -\mathbb{E} V(X_t)$$

$$\dot{\Sigma}_t = 2 I_d - \mathbb{E}\Big( \nabla V(X_t) \otimes (X_t - m_t) + (X_t - m_t) \otimes V(X_t) \Big)$$

$$N(m_t, \Sigma_t) ?$$

how to compute expectations ?

$$X_t \sim \pi_t \qquad Y_t \sim N(m_t, \Sigma_t)$$

$$\dot{m}_t = -\mathbb{E} \nabla V(Y_t)$$

$$\dot{\Sigma}_t = 2 I_d - \mathbb{E}\Big( \nabla V(Y_t) \otimes (Y_t - m_t) + (Y_t - m_t) \otimes V(Y_t) \Big)$$

$$P_t \sim N(m_t, \Sigma_t) \qquad (P2)$$

$$m_t \neq \mathbb{E}(X_t), \quad \Sigma_t \neq Cov(X_t)$$

Thm ( L-C-B-B-R '22):
The law $(p_t)_{t \geq 0}$ of $(p_2)$ is a
gradient flow of $KL(\cdot \| \gamma)$
on the Wasserstein space $(P_2(\mathbb{R}^d), W_2)$
which is constrained to the space
of Gaussians.

Bures-Wasserstein space:
$$BW(\mathbb{R}^d) - m \in \mathbb{R}^d$$
$$\Sigma \in S_{++} - \text{cone of}$$
$$\text{p.d. matrices}$$

Transport: $P_0 := P_{m_0, \Sigma_0} \dashrightarrow P_1 : P_{m_1, \Sigma_1}$

$$\nabla \varphi(x) = m_1 + \Sigma_0^{-1/2} \left( \Sigma_0^{1/2} \Sigma_t \Sigma_0^{1/2} \right) \Sigma_0^{-1/2} (x - m_0)$$

$$= \text{affine map}$$

$$W_2^2(P_0, P_1) = \| m_1 - m_2 \|^2 + tr \left( \Sigma_0 + \Sigma_1 - 2 \left( \Sigma_0^{1/2} \Sigma_1 \Sigma_0^{1/2} \right)^{1/2} \right)$$

Bures- Jko scheme :

(B) $P_{n, K\tau} := \underset{P \in BW(\mathbb{R}^d)}{\arg \min} \left[ KL \left( P \| \pi \right) + \frac{1}{2} W_2^2 \left( P, P_{n,k} \right) \right]$

Thm $(L-L-B-B-R' 22)$ :
Let $\pi \propto \exp(-v)$ be the target density on $\mathbb{R}^d$. The limiting curve $(\pi_t)_{t \geq 0}$ where $P_t = N(m_t, \Sigma_t)$ is obtained by (B).
BW gradient flow $(P_t)_{t \geq 0}$ of the KL divergence $KL(\cdot \| \pi)$ satisfies (P2).

___

Convergence properties.

$\mathcal{F} : BW(\mathbb{R}^d) \to \mathbb{R} \cup \{\infty\}$ and $\alpha \in \mathbb{R}$

$\mathcal{F}$ is $\alpha$-convex if for all constant speed geodesics $(P_t)_{t \geq 0}$ in $BW(\mathbb{R}^d)$

$\mathcal{F}(P_t) \leq (1-t) \mathcal{F}(P_0) + t \mathcal{F}(P_1) - \frac{\alpha t (1-t)}{2} W_2^2(P_0, P_1)$

$t \in [0, 1]$.

(Lem) For any $\alpha \in \mathbb{R}$ if $\nabla^2 V \succ \alpha I$
then $KL(\cdot \| \pi)$ is $\alpha$-convex on
$BW(\mathbb{R}^d)$ and there is a unique
soln. to $BW$-gradient flow of
$KL(\cdot \| \pi)$ starting at $P_0$ and

1. If $\alpha > 0$, $\forall t \geq 0$
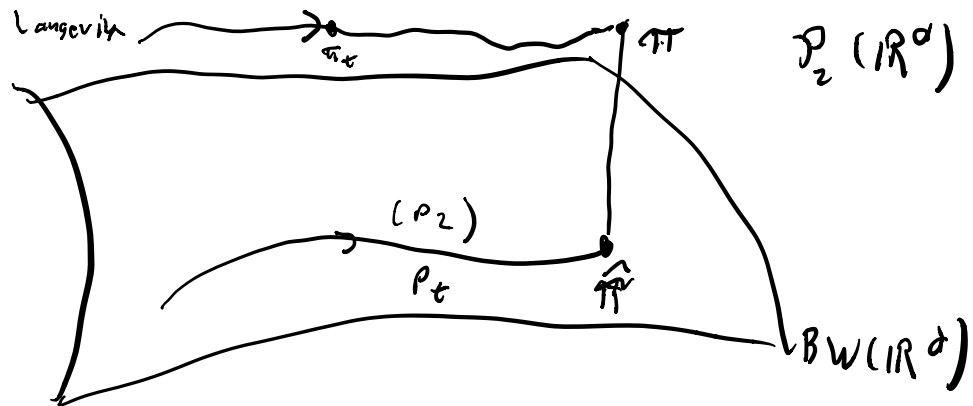$$W_2^2\left(P_t, \hat{\pi}\right) \leq \exp(-2\alpha t) W_2^2(P_0, \hat{\pi})$$

2. If $\alpha > 0$, $\forall t \geq 0$
$$KL(P_t \| \pi) \sim KL(\hat{P} \| \pi) \leq \exp(-2\alpha t) \left[ KL(P_0 \| \pi) - KL(\hat{\pi} \| \pi) \right]$$

3. If $\alpha = 0$, $\forall t > 0$
$$KL(P_t \| \pi) - KL(\hat{\pi} \| \pi) \leq \frac{1}{2t} W_2^2(P_0, \hat{\pi}).$$

$\mathcal{P}_2(\mathbb{R}^d)$

$BW(\mathbb{R}^d)$

JKO: store $\left(\pi_t\right)$

B-JKO: store $\left(m_t, \Sigma_t\right)$

Algorithms –

1) $\dot{m}_t = -\mathbb{E}\nabla V(Y_t)$

$\dot{\Sigma}_t = 2I_d - \mathbb{E}\left(\nabla V(Y_t)\otimes(Y_t - m_t) + (Y_t - m_t)\otimes\nabla V(Y_t)\right]$

Practically good

2) Provable algorithm - PA
$\alpha > 0$, $h > 0$, $m_0$, $\Sigma_0$

For $k = 1, \ldots, N$

$$\hat{X}_k \sim P_k$$

$$m_{k+1} \leftarrow m_k - h \nabla V(\hat{X}_k)$$

$$M_k \leftarrow I - h(\nabla^2 V(\hat{X}_k) - \Sigma_k^{-1})$$

$$\Sigma_k^+ \leftarrow M_k \Sigma_\alpha M_k$$

$$\Sigma_{k=} \leftarrow \text{clip}^{1/\alpha}(\Sigma_k^+)$$

$$\text{clip}^\tau(\Sigma) = \sum_{i=1}^{d} (\lambda_i \wedge \tau) u_i u_i^T$$

Thm: Assume $0 < \alpha I < \nabla^2 V < I$, $h < \frac{\alpha}{6}$,
initialize with $\frac{\alpha}{4} I < \Sigma_{m_0} < \frac{1}{\alpha} I$,
then $\forall \; k \in \mathbb{N}$

$$\mathbb{E} \, W_2^2(P_k, \hat{\pi}) \leq \exp(-\alpha k h) \, W_2^2(P_0, \hat{\pi}) + \frac{21 \, dh}{\alpha^2}.$$

$$\mathbb{E}\, W_2^2(\rho_K, \hat{\pi}) \leq \varepsilon^2 \quad \text{provided}$$

$$h \approx \frac{\alpha^2 \varepsilon^2}{d} \quad \& \quad K \gtrsim \frac{d}{\alpha^2 \varepsilon^2} \log\left(W_2(\rho_0, \hat{\pi})/\varepsilon\right)$$

$$\text{Complexity:} \quad O\left(\frac{d}{\alpha^2 \varepsilon^2} \log\left(\frac{W_2(\rho_0, \hat{\pi})}{\varepsilon}\right)\right)$$

---

## PA is SGD

Recall: $\left(\rho_t = \rho_{m_t, \Sigma_t}\right)_{t \geq 0}$

Bures-Wasserstein gradient:

$$g_\rho := \nabla_{BW} KL(\cdot \| \pi)$$

$$\simeq \left(\mathbb{E}_\rho \nabla V, \ \mathbb{E}_\rho \nabla^2 V - \text{cov}_\rho^{-1}\right)$$

$$\hat{g}_\rho := \left(\nabla V(\hat{x}), \ \nabla^2 V(\hat{x}) - \text{cov}_\rho^{-1}\right),$$

$$\tilde{x} \sim \rho$$

$$\rho_k^+ = \rho_{m_{k+1}, \Sigma_k^+}, \quad h \leq 1$$

$$P_k^+ = \exp_{P_k}(-h\,\hat{g}_k) \quad \ast$$

$$\hat{g}_k \in T_{P_k}\,BW(\mathbb{R}^d) \quad \text{is the stochastic gradient}$$

$$\hat{g}_k(x) = \nabla v(\hat{x}_k) + \left(\nabla^2 v(\hat{x}_k) - \Sigma_k^{-1}\right)$$

$$(x - m_k)$$

---

Question

$$X_t \sim \pi_t \qquad\qquad Y_t \sim GP(m(\cdot),\, \kappa(\cdot,\cdot))$$