

变分推理方法

最速梯度下降法

$$\text{minimize}_x f(x)$$

$$x_{k+1} = x_k - \alpha \nabla f(x_k)$$

$$\nabla f(x_k) = \arg \max_v \frac{\lim_{\epsilon \rightarrow 0} \frac{f(x_k + \epsilon v) - f(x_k)}{\epsilon}}{\sqrt{\langle v, v \rangle}}$$

$$= \arg \max_v \frac{\frac{\delta f}{\delta x}(x_k) \cdot v}{\sqrt{\langle v, v \rangle}}$$

变化最大的方向

$$\frac{\delta f}{\delta x}(x_k) \cdot v \leq \sqrt{\langle v, v \rangle} \sqrt{\frac{\delta f}{\delta x}(x_k) \cdot \frac{\delta f}{\delta x}(x_k)}$$

$$\text{当且仅当 } v = \frac{\delta f}{\delta x}(x_k) = \nabla f(x_k)$$

度量

切空间: $T_{x_k} M$ 局部线性近似

曲线: $x_t: [0, 1] \rightarrow M$

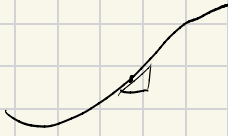
$$\dot{x}_t = v_t$$

曲线长度: $ds^2 = g_x(v, v)$

$$= \langle M_x v, v \rangle$$

练习 1)

$$\dot{x}_t = v_t$$



曲线 $x(t)$, 长度 $\int_0^1 \|\dot{x}(t)\|_2 dt$

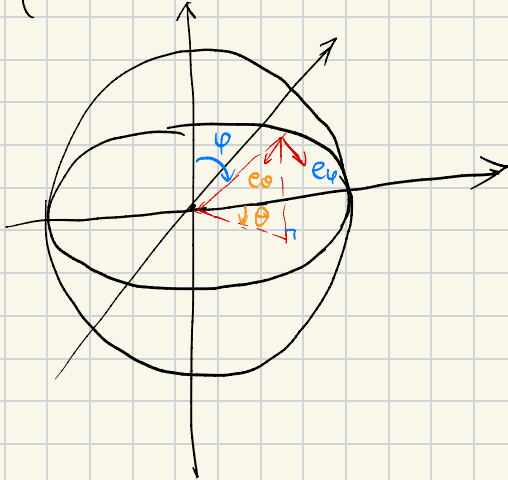
$$= \int_0^1 \sqrt{v_t \cdot v_t} dt$$

$$= \int_0^1 \sqrt{g_x(v_t, v_t)} dt$$

练习 2) 球面

$$(\varphi, \theta) \rightarrow (\sin \varphi \cos \theta, \sin \varphi \sin \theta, \cos \varphi)$$

切空间



$$(\dot{\varphi}, \dot{\theta}) \rightarrow \dot{\varphi} e_\varphi + \dot{\theta} e_\theta$$

$$e_\varphi = (\cos \varphi \cos \theta, \cos \varphi \sin \theta, -\sin \varphi)$$

$$e_\theta = (-\sin \varphi \sin \theta, \sin \varphi \cos \theta, 0)$$

曲线 $\dot{x}_t = v_t = (\dot{\varphi}_t \dot{\theta}_t)$

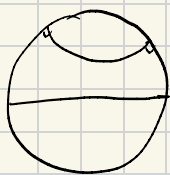
欧式空间 距离

$$g_x(v, v) = \langle \dot{\varphi} e_\varphi + \dot{\theta} e_\theta, \dot{\varphi} e_\varphi + \dot{\theta} e_\theta \rangle$$
$$= \dot{\varphi}^2 + \dot{\theta}^2 \sin^2 \varphi$$

$$= \langle M_x v, v \rangle$$

$$M_x = \begin{bmatrix} 1 & 0 \\ 0 & \sin^2 \varphi \end{bmatrix}$$

练习三：庞加莱圆盘 (双曲几何)



在地上无限大

内积 $g_x(v, v) = \langle M_x v, v \rangle$

$$M_x = \begin{bmatrix} \frac{4}{(1-x^2-y^2)^2} & \\ & \frac{4}{(1-x^2-y^2)^2} \end{bmatrix}$$

$$(\dot{r}, \dot{\theta}) \rightarrow \dot{r} e_r + \dot{\theta} e_\theta$$

$$e_r = (\cos \theta, \sin \theta) \quad e_\theta = (-r \sin \theta, r \cos \theta)$$

$$g_x(v, v) = (\dot{r} \cos \theta - \dot{\theta} r \sin \theta)^2 \frac{4}{(1-r^2)^2}$$
$$+ (\dot{r} \sin \theta + \dot{\theta} r \cos \theta)^2 \frac{4}{(1-r^2)^2}$$
$$= \dot{r}^2 \frac{4}{(1-r^2)^2} + \dot{\theta}^2 r^2 \frac{4}{(1-r^2)^2}$$

$$= \langle M_x v, v \rangle$$

$$M_x = \begin{bmatrix} \frac{4}{(1-r^2)^2} & \\ & \frac{4r^2}{(1-r^2)^2} \end{bmatrix}$$

流形上的最速梯度下降法

$$\frac{\partial f}{\partial x}(x_k) \cdot v \leq \sqrt{\langle M(x)v, v \rangle} \cdot \sqrt{\langle M(x)^{-1} \frac{\partial f}{\partial x}, \frac{\partial f}{\partial x} \rangle}$$

$$\text{当且仅当 } M(x)v = \frac{\partial f}{\partial x}(x_k)$$

能量泛函

$$\frac{\delta E}{\delta \rho} b = \lim_{\varepsilon \rightarrow 0} \frac{\int (P + \varepsilon b) \log\left(\frac{P + \varepsilon b}{\rho^*}\right) - \int P \log\left(\frac{P}{\rho^*}\right)}{\varepsilon}$$

$$= \lim_{\varepsilon \rightarrow 0} \frac{\int P \log\left(\frac{P + \varepsilon b}{P}\right) + \varepsilon b \log\left(\frac{P + \varepsilon b}{\rho^*}\right)}{\varepsilon}$$

$$= \lim_{\varepsilon \rightarrow 0} \frac{\int P \left(\frac{\varepsilon b}{P}\right) + \varepsilon b \log\left(\frac{P + \varepsilon b}{\rho^*}\right)}{\varepsilon}$$

$$= \int b \log \frac{P}{\rho^*} d\sigma$$

Wasserstein-2 度量 (P compactly support)

Otto calculus, Wasserstein-2 度量

考虑: $\inf_v \int \frac{\|v\|_2^2}{2} p \, d\theta$, $b = -\nabla_\theta \cdot (p v)$

$$\inf_v \sup_\psi \int \frac{\|v\|_2^2}{2} p \, d\theta + \int (b + \nabla_\theta \cdot (p v)) \psi \, d\theta$$
$$= \inf_v \sup_\psi \int \frac{\|v\|_2^2}{2} p \, d\theta + \int b \psi \, d\theta - \int p v \cdot \nabla_\theta \psi \, d\theta$$

$$= \sup_\psi \inf_v \int p \left(\frac{v \cdot v}{2} - v \nabla_\theta \psi \right) d\theta + \int b \psi \, d\theta$$

$$= \sup_\psi \int b \psi - p \frac{\|\nabla_\theta \psi\|^2}{2} \, d\theta \quad \begin{cases} v = \nabla_\theta \psi \\ b = -\nabla_\theta \cdot (p \nabla_\theta \psi) \end{cases}$$

由于 $b \in T_p P$ $b = -\nabla_\theta \cdot (p \nabla_\theta \psi)$, $\nabla_\theta \psi$ 有唯一解

$$= \int -\nabla_\theta \cdot (p \nabla_\theta \psi) \psi - p \frac{\|\nabla_\theta \psi\|^2}{2} \, d\theta$$

$$= \int \frac{\|\nabla_\theta \psi\|_2^2}{2} p \, d\theta$$

$$\text{因此 } g_e^{W_2}(b, b) = \int \|\nabla_\theta \psi\|_2^2 p \, d\theta \quad b = -\nabla_\theta \cdot (p \nabla_\theta \psi)$$

$$= \int \underline{-\nabla_\theta \cdot (p \nabla_\theta \psi)} \cdot \psi \, d\theta$$

3

$$= \langle \psi, \beta \rangle := \langle M_p^{W_2} \beta, \beta \rangle$$

$$M(p)^{W_2} \beta = \psi \quad \beta = -\nabla_{\theta} (p \nabla_{\theta} \psi)$$

$$M(p)^{W_2} \psi = \beta = -\nabla_{\theta} (p \nabla_{\theta} \psi)$$

Wasserstein 梯度流

$$\begin{aligned} \frac{\partial p_t}{\partial t} &= \nabla_{\theta} (p_t \nabla_{\theta} (\log p_t - \log p^*)) \\ &= \nabla_{\theta} (p_t \left(\frac{\nabla_{\theta} p_t}{p_t} - \nabla_{\theta} [-\Phi_R] \right)) \\ &= \nabla_{\theta} (\nabla_{\theta} p_t + p_t \nabla_{\theta} \Phi_R) \end{aligned}$$

高斯近似的 Wasserstein 梯度流

$$\int \nabla_{\theta} \cdot f \cdot M = \sum \int \partial_i f_i M = \sum \int f_i \partial_i M$$

$$= \int \underbrace{(p_{at} \nabla_{\theta} \Phi_R + \nabla_{\theta} p_{at})}_i [e_i \cdot (\theta - m_t)^T + (\theta - m_t) e_i^T]$$

$$= \int (p_{at} \nabla_{\theta} \Phi_R + \nabla_{\theta} p_{at}) (\theta - m_t)^T + (\theta - m_t) (p_t \nabla_{\theta} \Phi_R + \nabla_{\theta} p_t)^T$$

$$\text{由于 } \int \nabla_{\theta} p_{at} (\theta - m_t)^T = \int p_{at} C_t^{-1} (\theta - m_t) (\theta - m_t)^T = I$$

$$\begin{aligned} \int p_{at} \nabla_{\theta} \Phi_R (\theta - m_t)^T &= \int \nabla_{\theta} \Phi_R (C_t \nabla_{\theta} p_{at})^T \\ &= -\int \nabla_{\theta} \nabla_{\theta} \Phi_R p_{at} d\theta C_t \end{aligned}$$

收敛性:

$$\partial_t \text{KL}[p_t \| p^*] = \frac{\partial}{\partial t} \int p_t \log \frac{p_t}{p^*} d\theta$$

$$= \int \dot{p}_t \log \frac{p_t}{p^*} + \underbrace{p_t \frac{1}{p_t} \dot{p}_t}_{\dot{p}_t = 0} d\theta$$

$$= \int \nabla_{\theta} \cdot (p_t \nabla_{\theta} \log p_t - \log p^*) \log \left(\frac{p_t}{p^*} \right) d\theta$$

$$= - \int p_t \|\nabla_{\theta} \log \frac{p_t}{p^*}\|_2^2 d\theta$$

log-sobolev

$$\leq -2\alpha \text{KL}[p_t \| p^*]$$

Gronwall 引理

$$\partial_t (e^{2\alpha t} \text{KL}[p_t \| p^*]) \leq 0$$

$$e^{2\alpha t} \text{KL}[p_t \| p^*] \leq \text{KL}[p_0 \| p^*]$$

Fisher - Rao metric

$$\int_{\tilde{\theta} \in A} T_{\#} p(\tilde{\theta}) d\tilde{\theta} = \int_{T\theta \in A} p(\theta) d\theta$$
$$= \int_{T\theta \in A} T_{\#} p(T\theta) |T'(\theta)| d\theta$$

$$p(\theta) = T_{\#} p(T\theta) |T'(\theta)|$$

$$z(\theta) = T_{\#} z(T\theta) |T'(\theta)|$$

$$\int \frac{T_{\#} z(\tilde{\theta}) \cdot T_{\#} z(\tilde{\theta})}{T_{\#} p(\tilde{\theta})} d\tilde{\theta} \quad \tilde{\theta} = T\theta$$

$$= \int \frac{z(\theta) / |T'(\theta)| \cdot z(\theta) / |T'(\theta)|}{p(\theta) / |T'(\theta)|} |T'(\theta)| d\theta$$

$$= \int \frac{z(\theta) \cdot z(\theta)}{p(\theta)} d\theta$$

对偶梯度

$$\begin{aligned} & \frac{d}{dt} \left(\int p_t \log \frac{p_t}{p^*} + \int p^* \log \frac{p^*}{p_t} \right) \\ &= \int \dot{p}_t \log \frac{p_t}{p^*} - \int p^* \frac{1}{p_t} \dot{p}_t \\ &= - \int p_t \left(\log \frac{p_t}{p^*} \right)^2 - \int p_t \log \frac{p_t}{p^*} \cdot c - \int p^* \log \frac{p^*}{p_t} + \int p^* c \\ &= \underbrace{- \int p_t \left(\log \frac{p_t}{p^*} \right)^2 + \left(\int p_t \log \frac{p_t}{p^*} \right)^2}_{\leq 0} - \text{KL}[p^* || p_t] - \text{KL}[p_t || p^*] \\ &\leq -\text{KL}[p^* || p_t] - \text{KL}[p_t || p^*] \end{aligned}$$

生死过程

θ 以 $1 - e^{-|\lambda_t| dt}$ 的概率复制 / 湮灭

当 $\lambda_t > 0$ / $\lambda_t < 0$.

$$\begin{aligned} & \int_{\theta \in A} p_{t+\lambda_t}(\theta) d\theta - \int_{\theta \in A} p_t(\theta) d\theta \\ &= \int_{\substack{\lambda_t(\theta) > 0 \\ \theta \in A}} p_t(\theta) (1 - e^{-|\lambda_t| dt}) - \int_{\substack{\lambda_t(\theta) < 0 \\ \theta \in A}} p_t(\theta) (1 - e^{-|\lambda_t| dt}) \\ &= \int p_t(\theta) \text{sign}(\lambda_t(\theta)) (|\lambda_t| dt + o(dt)) \end{aligned}$$

$$\frac{\partial \rho_t}{\partial t} = \rho_t \delta_t$$

参数化变分推理

$$\int \rho_{\theta} [(-\Phi_R - \log \rho_{\theta}) - \underbrace{E_{\rho_{\theta}}[-\Phi_R - \log \rho_{\theta}]}] \theta \, d\theta$$

$$= \int \rho_{\theta} [-\Phi_R - \log \rho_{\theta}] (\theta - m_t) \, d\theta$$

$$= -C_t \int \nabla_{\theta} \rho_{\theta} (-\Phi_R - \log \rho_{\theta}) \, d\theta$$

$$= -C_t \int \rho_{\theta} (\nabla_{\theta} \Phi_R + \nabla_{\theta} \log \rho_{\theta}) \, d\theta$$

$$= -C_t E \nabla_{\theta} \Phi_R$$

$$\int \rho_{\theta} [(-\Phi_R - \log \rho_{\theta}) - E_{\rho_{\theta}}[-\Phi_R - \log \rho_{\theta}]] (\theta - m_t)(\theta - m_t)^T \, d\theta$$

$$= \int \rho_{\theta} [-\Phi_R - \log \rho_{\theta}] [(\theta - m_t)(\theta - m_t)^T - C_t] \, d\theta$$

$$= C_t \int [-\Phi_R - \log \rho_{\theta}] \nabla_{\theta} \nabla_{\theta} \rho_{\theta} \, d\theta C_t$$

$$= C_t \int \nabla_{\theta} \nabla_{\theta} [-\Phi_R - \log \rho_{\theta}] \rho_{\theta} \, d\theta C_t$$

$$= C_t - C_t E \nabla_{\theta} \nabla_{\theta} \Phi_R C_t$$

当 Φ_R 强凸, KL 散度在 W_2 下是强凸,

Gaussian space, W_2 , BW_2 , KL 散度也是强凸

因为 G_A \curvearrowright G_B
测地线全是高斯

自然梯度下降法

$$P = \{ p : \|p\| = 1 \} \quad T_p P = \{ v : \int v = 0 \}$$

$$P = \{ a \in \mathbb{R}^{N_a} \} \quad T_{p_a} P = \{ v : \in \mathbb{R}^{N_a} \}$$

$$p_a \rightarrow p_{a+\varepsilon v} = p_a + \varepsilon v$$

$$\delta = \lim_{\varepsilon \rightarrow 0} \frac{p_{a+\varepsilon v} - p_a}{\varepsilon} = \nabla_a p_a \cdot v$$

$$g_{p_a}^{\text{FR}} (\nabla_a p_a \cdot v_1, \nabla_a p_a \cdot v_2)$$

$$= \int \frac{(\nabla_a p_a \cdot v_1)^T \cdot (\nabla_a p_a \cdot v_2)}{p_a} d\theta$$

$$= v_1^T \int \frac{\nabla_a p_a \nabla_a p_a}{p_a} d\theta v_2$$

$$KL[P_{a+da} \parallel P_a] = \int P_{a+da} \log \frac{P_{a+da}}{P_a}$$

$$P_{a+da} = P_a + \nabla_a P_a \cdot da + \frac{1}{2} da^T \nabla_a^2 P_a da$$

$$\log \frac{P_{a+da}}{P_a} = \frac{\nabla_a P_a \cdot da}{P_a} + \frac{\frac{1}{2} da^T \nabla_a^2 P_a da}{P_a} - \frac{1}{2} \left(\frac{\nabla_a P_a \cdot da}{P_a} \right)^2$$

$$\log 1+x = x - \frac{x^2}{2}$$

↑↑

$$= \int \nabla_a P_a \cdot da + \frac{1}{2} da^T \nabla_a^2 P_a da + \frac{(\nabla_a P_a \cdot da)^2}{P_a} d\theta - \frac{1}{2} \left(\frac{\nabla_a P_a \cdot da}{P_a} \right)^2 P_a$$

$$= \frac{1}{2} da^T \int \frac{\nabla_a P_a \cdot \nabla_a P_a}{P_a} d\theta da$$

1. 高斯变分推理

对于高斯分布

$$\rho_a = \frac{1}{(2\pi)^{N_\theta/2} \sqrt{|C|}} \exp\left\{-\frac{1}{2}(\theta - m)^T C^{-1}(\theta - m)\right\}$$

我们有

$$\begin{aligned}\nabla_m \rho_a &= \rho_a C^{-1}(\theta - m) = -\nabla_\theta \rho_a \\ \nabla_C \rho_a &= \rho_a \left(-\frac{1}{2} C^{-1} + \frac{1}{2} C^{-1}(\theta - m)(\theta - m)^T C^{-1} \right) = \frac{1}{2} \nabla_\theta \nabla_\theta \rho_a\end{aligned}$$

我们用了

$$\nabla_C |C| = |C| C^{-T} \quad \nabla_C x^T C^{-1} x = -C^{-1} x x^T C^{-1}$$

对于高斯变分推理，我们有

$$\begin{aligned}
\nabla_m KL[\rho_a \|\rho^*] &= \int \nabla_m \rho_a (\log(\rho_a) + \Phi_R) d\theta \\
&= - \int \nabla_\theta \rho_a (\log(\rho_a) + \Phi_R) d\theta \\
&= \int \rho_a \nabla_\theta (\log(\rho_a) + \Phi_R) d\theta \\
&= \int \nabla_\theta \rho_a + \rho_a \nabla_\theta \Phi_R d\theta \\
&= \int \rho_a \nabla_\theta \Phi_R d\theta \\
&= \mathbb{E}_{\rho_a} [\nabla_\theta \Phi_R] \\
\nabla_C KL[\rho_a \|\rho^*] &= \int \nabla_C \rho_a (\log(\rho_a) + \Phi_R) d\theta \\
&= \int \frac{1}{2} \nabla_\theta \nabla_\theta \rho_a (\log(\rho_a) + \Phi_R) d\theta \\
&= \int \frac{1}{2} \rho_a \nabla_\theta \nabla_\theta (\log(\rho_a) + \Phi_R) d\theta \\
&= \int \frac{1}{2} \rho_a (-C) + \frac{1}{2} \rho_a \nabla_\theta \nabla_\theta \Phi_R d\theta \\
&= -\frac{1}{2} C + \frac{1}{2} \mathbb{E}_{\rho_a} [\nabla_\theta \nabla_\theta \Phi_R]
\end{aligned}$$

2. 自然梯度下降方法

我们首先推导高斯的Fisher信息矩阵 我们有

$$\int \frac{\nabla_m \rho_a \nabla_m \rho_a^T}{\rho_a} d\theta = \int \rho_a C^{-1}(\theta - m)(\theta - m)^T C^{-1} d\theta d\theta = C^{-1}$$

$$\int \frac{\nabla_m \rho_a \nabla_C \rho_a^T}{\rho_a} d\theta = \int \rho_a C^{-1}(\theta - m) \otimes \left(-\frac{1}{2} C^{-1} - \frac{1}{2} C^{-1}(\theta - m)(\theta - m)^T C^{-1} \right) d\theta = 0$$

对于协方差项，我们定义

$$X := \frac{1}{4} \int \rho_a \left(C^{-1}(\theta - m)(\theta - m)^T C^{-1} - C^{-1} \right) \otimes \left(C^{-1}(\theta - m)(\theta - m)^T C^{-1} - C^{-1} \right) d\theta$$

$$= \frac{1}{4} \int \mathcal{N}(y; 0, I) C^{-1/2} (yy^T - I) C^{-1/2} \otimes C^{-1/2} (yy^T - I) C^{-1/2} dy, \text{ where } y = C^{-1/2}(\theta - m).$$

它满足

$$X[ij, lm] = \frac{1}{4} \sum_{r,s,p,q} C^{-1/2}[i, r] C^{-1/2}[j, s] C^{-1/2}[l, p] C^{-1/2}[m, q] \int (y_r y_s - \delta_{r,s})(y_p y_q - \delta_{p,q}) \mathcal{N}(y; 0, I) dy$$

$$= \frac{1}{4} \sum_{r,s,p,q} C^{-1/2}[i, r] C^{-1/2}[j, s] C^{-1/2}[l, p] C^{-1/2}[m, q] (\delta_{r,p} \delta_{s,q} + \delta_{r,q} \delta_{s,p})$$

$$= \frac{1}{4} \left(C^{-1}[i, l] C^{-1}[j, m] + C^{-1}[i, m] C^{-1}[j, l] \right).$$

因此

$$(XY)_{ij} = \sum_{l,m} X[ij, lm] Y[l, m]$$

$$= \frac{1}{4} \sum_{l,m} \left(C^{-1}[i, l] Y[l, m] C^{-1}[j, m] + C^{-1}[i, m] Y[l, m] C^{-1}[j, l] \right)$$

$$= \frac{1}{4} (C^{-1} Y C^{-1} + C^{-1} Y^T C^{-1})_{ij}.$$

我们用了

由于

$$\frac{dC_t^{-1}}{dt} = -C_t^{-1} \frac{dC_t}{dt} C_t^{-1}$$

自然梯度下降法可以写成

$$\begin{aligned} \frac{dm_t}{dt} &= C_t C^{*-1} (m^* - m_t), \\ \frac{dC_t^{-1}}{dt} &= -C_t^{-1} + C^{*-1}. \end{aligned}$$

对于协方差，我们有

$$C_t^{-1} = (1 - e^{-t})C^{*-1} + e^{-t}C_0^{-1}.$$

对于期望，我们有

$$m^* - m_t = e^{-t}C_t C_0^{-1} (m^* - m_0).$$

计算它的导数，我们能得到

$$\begin{aligned} \frac{dm_t}{dt} &= -e^{-t}C_t C_0^{-1} (m_0 - m^*) + e^{-t} \frac{dC_t}{dt} C_0^{-1} (m_0 - m^*) \\ &= -e^{-t}C_t C_0^{-1} (m_0 - m^*) + e^{-t} (C_t - C_t C^{*-1} C_t) C_0^{-1} (m_0 - m^*) \\ &= C_t C^{*-1} e^{-t} C_t C_0^{-1} (m^* - m_0) \\ &= C_t C^{*-1} (m^* - m_t). \end{aligned}$$

```
In [1]: using LinearAlgebra
using PyPlot
function expectation(method_type::String, m_oo, C_oo, m, C)
    return C_oo \ (m - m_oo), inv(C_oo)
    if method_type == "gradient_descent"
```

```

    # under-determined case
    return -C_oo*(m - m_oo), 1/2.0*(inv(C) - inv(C_oo))

elseif method_type == "natural_gradient_descent"
    # over-determined case
    return -(m - m_oo), C - C*(C_oo\C)

elseif method_type == "natural_gradient_descent_Cinv"
    # over-determined case
    return -(m - m_oo), -inv(C) + C_oo

elseif method_type == "wasserstein_gradient_descent"
    # over-determined case
    return -C_oo*(m - m_oo), 2I - C_oo\C - C/C_oo

else
    error("Problem type : ", problem_type, " has not implemented!")
end
end

function Continuous_Dynamics(method_type::String, m_oo, C_oo, m_0, C_0, Δt, N_t)

    N_θ = length(m_0)
    m = zeros(N_t+1, N_θ)
    C = zeros(N_t+1, N_θ, N_θ)

    m[1, :] = m_0
    C[1, :, :] = C_0

    for i = 1:N_t
        EVΦ, EV∇Φ = expectation(method_type, m_oo, C_oo, m[i, :], C[i, :, :])
        if method_type == "gradient_descent"
            m[i+1, :] = m[i, :] - EVΦ * Δt
            C[i+1, :, :] = C[i, :, :] + 1/2.0*(inv(C[i, :, :]) - EV∇Φ) * Δt
        elseif method_type == "natural_gradient_descent"
            m[i+1, :] = m[i, :] - C[i, :, :]*EVΦ * Δt
            C[i+1, :, :] = inv( inv(C[i, :, :]) + (EV∇Φ - inv(C[i, :, :])) * Δt)
        else

```

```

        error("Problem type : ", problem_type, " has not implemented!")
    end
end

return m, C
end

Δt, N_t = 5e-1, 30
fig, ax = PyPlot.subplots(ncols=2, sharex=true, sharey=true, figsize=(6,3))

m_oo = [1; 1]
C_oo = [100.0 0; 0 1.0]

m_0 = [0.0; 0.0]
C_0 = [1.0 0; 0 1.0]

m_gd, C_gd = Continuous_Dynamics("gradient_descent", m_oo, C_oo, m_0, C_0, Δt, N_t)
m_ngd, C_ngd = Continuous_Dynamics("natural_gradient_descent", m_oo, C_oo, m_0, C_0, Δt, N_t)

e_gd, e_ngd = zeros(N_t+1, 2), zeros(N_t+1, 2)

for i = 1:N_t+1

    e_gd[i, 1] = ( norm(m_gd[i,:] - m_oo) / norm(m_oo) )
    e_gd[i, 2] = ( norm(C_gd[i,:,:] - C_oo) / norm(C_oo) )

    e_ngd[i, 1] = ( norm(m_ngd[i,:] - m_oo) / norm(m_oo) )
    e_ngd[i, 2] = ( norm(C_ngd[i,:,:] - C_oo) / norm(C_oo) )
end

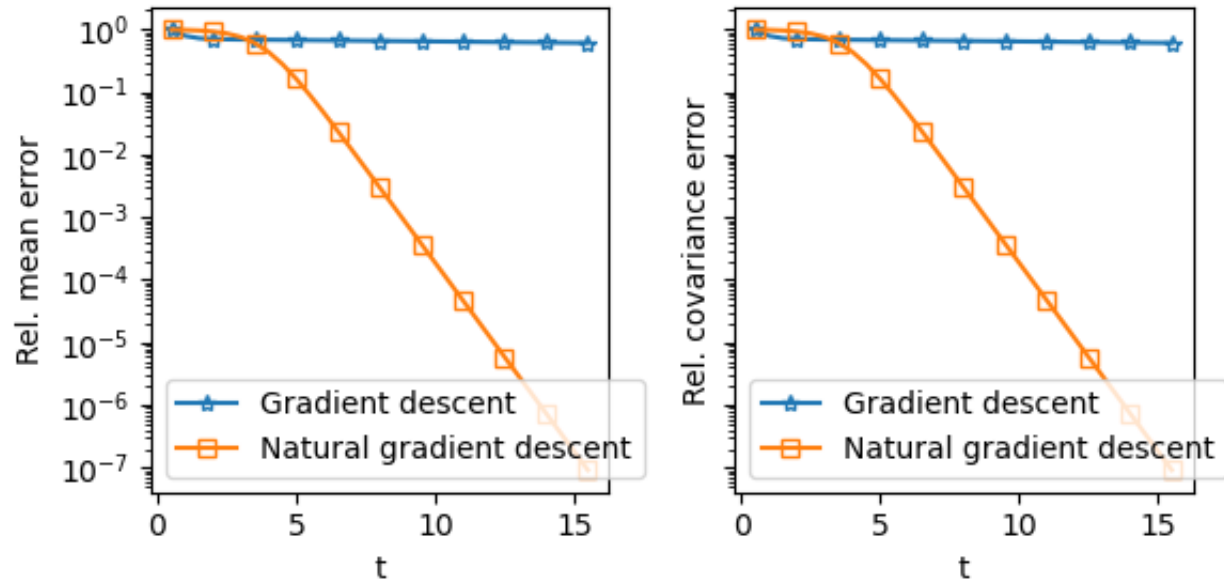
ts = Array(1:N_t+1)*Δt
ax[1].semilogy(ts, e_gd[:, 1], "--*", fillstyle="none", markevery=div(N_t, 10), label="Gradient descent")
ax[1].semilogy(ts, e_ngd[:, 2], "-s", fillstyle="none", markevery=div(N_t, 10), label="Natural gradient des")
ax[2].semilogy(ts, e_gd[:, 1], "--*", fillstyle="none", markevery=div(N_t, 10), label="Gradient descent")
ax[2].semilogy(ts, e_ngd[:, 2], "-s", fillstyle="none", markevery=div(N_t, 10), label="Natural gradient des")

```

```

ax[1].legend()
ax[2].legend()
ax[1].set_ylabel("Rel. mean error")
ax[2].set_ylabel("Rel. covariance error")
ax[1].set_xlabel("t")
ax[2].set_xlabel("t")
fig.tight_layout()
fig.savefig("Gaussian-VI.pdf")

```



3. 指数分布族

对于指数分布族

$$\rho(\theta; a) = h(\theta) \exp\{T(\theta) \cdot a - A(a)\}$$

我们有:

$$0 = \int \nabla_a \rho(\theta; a) d\theta = \mathbb{E}_\rho[T(\theta) - \nabla_a A(a)] \quad (1)$$

其次, Fisher信息矩阵 FIM 满足

$$\begin{aligned} FIM(a) &= \mathbb{E}_\rho[\nabla_a \log \rho(\theta; a)^T \nabla_a \log \rho(\theta; a)] \\ &= \int \nabla_a \rho(\theta; a) [T(\theta) - \nabla_a A(a)] d\theta \\ &= \int \rho(\theta; a) \nabla_a \nabla_a A(a) d\theta \quad \text{Using } \nabla_a \int \rho(\theta; a) [T(\theta) - \nabla_a A(a)] d\theta = 0 \\ &= \nabla_a \nabla_a A(a) \end{aligned} \quad (2)$$

4. Langevin 动力系统

考虑随机动力系统

$$d\theta_t = f(\theta_t, t)dt + \sqrt{g}dW_t$$

它对应的描述 $\theta_t \sim \rho_t$ 的 Fokker-Planck 方程是

$$\partial_t \rho_t = -\nabla \cdot (f(x, t)\rho_t) + \frac{1}{2} \nabla \cdot \nabla \cdot (g\rho_t)$$

使用 Itô's 公式, 我们有

$$\begin{aligned}
dh(\theta_t) &= \nabla h(\theta_t)^T (f(\theta_t, t)dt + \sqrt{g}dW_t) + \frac{1}{2}g\nabla^2 h(\theta_t)dt \\
\partial_t \mathbb{E}[h(\theta_t)] &= \mathbb{E}\left[f(\theta_t, t)\nabla h(\theta_t) + \frac{g}{2}\nabla^2 h(\theta_t)\right] \\
\int h(\theta)\partial_t \rho_t &= \int \rho_t f(\theta, t)\nabla h(\theta) + \frac{g}{2}\rho_t \nabla^2 h(\theta)d\theta \\
&= \int h(\theta)\nabla(-\rho_t f(\theta, t)) + h(\theta)\frac{1}{2}\nabla \cdot \nabla \cdot (g\rho_t)d\theta
\end{aligned}$$

我们可以取

$$f = -\nabla_{\theta}\Phi_R \quad g = 2$$

5. 仿射不变性

给定可逆仿射变换 $\tilde{\theta} = \mathcal{T}\theta = A\theta + b$,

$$\begin{aligned}
\tilde{\rho}(\tilde{\theta}) &= \rho(\theta)|A^{-1}| & \tilde{\rho}_t(\tilde{\theta}) &= \rho_t(\theta)|A^{-1}| & \tilde{C}_t &= AC_tA^T \\
\nabla_{\tilde{\theta}}f(\tilde{\theta}) &= A^{-T}\nabla_{\theta}f(\theta) & \nabla_{\tilde{\theta}} \cdot v(\tilde{\theta}) &= \nabla_{\theta} \cdot (A^{-1}v(\theta))
\end{aligned}$$

带入Kalman-Wasserstein梯度流

$$\begin{aligned}
\frac{\partial \tilde{\rho}_t}{\partial t} &= \nabla_{\tilde{\theta}} \cdot \left[\tilde{\rho}_t \tilde{C}_t (\nabla_{\tilde{\theta}} \log \tilde{\rho}^* - \nabla_{\tilde{\theta}} \log \tilde{\rho}_t) \right] \\
\iff \frac{\partial \rho_t}{\partial t} |A^{-1}| &= \nabla_{\tilde{\theta}} \cdot \left[\rho_t AC_t A^T (\nabla_{\tilde{\theta}} \log \tilde{\rho}^* - \nabla_{\tilde{\theta}} \log \tilde{\rho}_t) \right] |A^{-1}| \\
\iff \frac{\partial \rho_t}{\partial t} |A^{-1}| &= \nabla_{\theta} \cdot \left[\rho_t C_t (\nabla_{\theta} \log \tilde{\rho}^* - \nabla_{\theta} \log \tilde{\rho}_t) \right] |A^{-1}| \\
\iff \frac{\partial \rho_t}{\partial t} &= \nabla_{\theta} \cdot \left[\rho_t C_t (\nabla_{\theta} \log \rho^* - \nabla_{\theta} \log \rho_t) \right]
\end{aligned}$$

6. 练习

In [2]: `using` PyPlot, Random, Statistics, Distributions, ForwardDiff, LinearAlgebra

```
Random.seed!(42)
```

```
function Phi_R_Convex( $\theta$ )
```

```
     $\epsilon$  = 0.1
```

```
     $\theta_1, \theta_2$  =  $\theta$ 
```

```
    return ( ( $\epsilon * \theta_1 - \theta_2$ )2 +  $\theta_2$ 4 ) / 20
```

```
end
```

```
 $\Phi_R(\theta)$  = Phi_R_Convex( $\theta$ )
```

```
function  $\nabla\Phi_R(\theta)$ 
```

```
    return  $\Phi_R(\theta)$ , ForwardDiff.gradient( $\Phi_R$ ,  $\theta$ )
```

```
end
```

```
function update_ensemble( $\theta$ ,  $\Delta t$ , func_neg_log_rho::Function, preconditioner::Bool)
```

```
    N_ $\theta$ , N_ens = size( $\theta$ )
```

```
     $\Phi_R$ ,  $\nabla\Phi_R$  = zeros(N_ens), zeros(N_ $\theta$ , N_ens)
```

```
    for i = 1:N_ens
```

```
         $\Phi_R[i]$ ,  $\nabla\Phi_R[:,i]$  = func_neg_log_rho( $\theta[:, i]$ )
```

```
    end
```

```
     $\theta_{\text{mean}}$  = mean( $\theta$ , dims=2)
```

```
     $\theta\theta_{\text{cov}}$  = (( $\theta$  .-  $\theta_{\text{mean}}$ ) * ( $\theta$  .-  $\theta_{\text{mean}}$ )') / N_ens
```

```
    Prec = (preconditioner ?  $\theta\theta_{\text{cov}}$  : 1)
```

```
     $\sigma$  = (preconditioner ? cholesky(Hermitian( $\theta\theta_{\text{cov}}$ )).L : 1)
```

```
    noise = rand(Normal(0, 1), (N_ $\theta$ , N_ens))
```

```
     $\theta$  =  $\theta - \Delta t * \text{Prec} * \nabla\Phi_R + \text{sqrt}(2 * \Delta t) * (\sigma * \text{noise})$ 
```

```
return  $\theta$ 
```

```
end
```

```
Out[2]: update_ensemble (generic function with 1 method)
```

```
In [3]: N_θ = 2
m_0 = [10.0; 10.0]
σ0 = 2
C_0 = [σ0^2  0.0; 0.0  σ0^2]

Lx , Ux = -100.0, 100.0
Ly , Uy = -5.0, 5.0
N = 500
X = zeros(N, N)
Y = zeros(N, N)
ρ = zeros(N, N)
for i = 1:N
    for j = 1:N
        X[i, j], Y[i, j] = Lx + (Ux - Lx) * (i-1)/(N-1), Ly + (Uy - Ly) * (j-1)/(N-1)
        ρ[i, j] = Phi_R_Convex([X[i, j]; Y[i, j]])
    end
end
Z = sum(ρ)
ρ .= exp.(-ρ)/Z

fig, ax = PyPlot.subplots(ncols=2, nrows=1, sharex="col", sharey=true, figsize=(10,3))
ax[1].contour(X, Y, ρ, 10, colors="grey")
ax[2].contour(X, Y, ρ, 10, colors="grey")

N_ens = 100
θ0 = Array(rand(MvNormal(m_0, C_0), N_ens))
θ = zeros(N_θ, N_ens)
Δt = 0.01
N_t = 2000
```

```

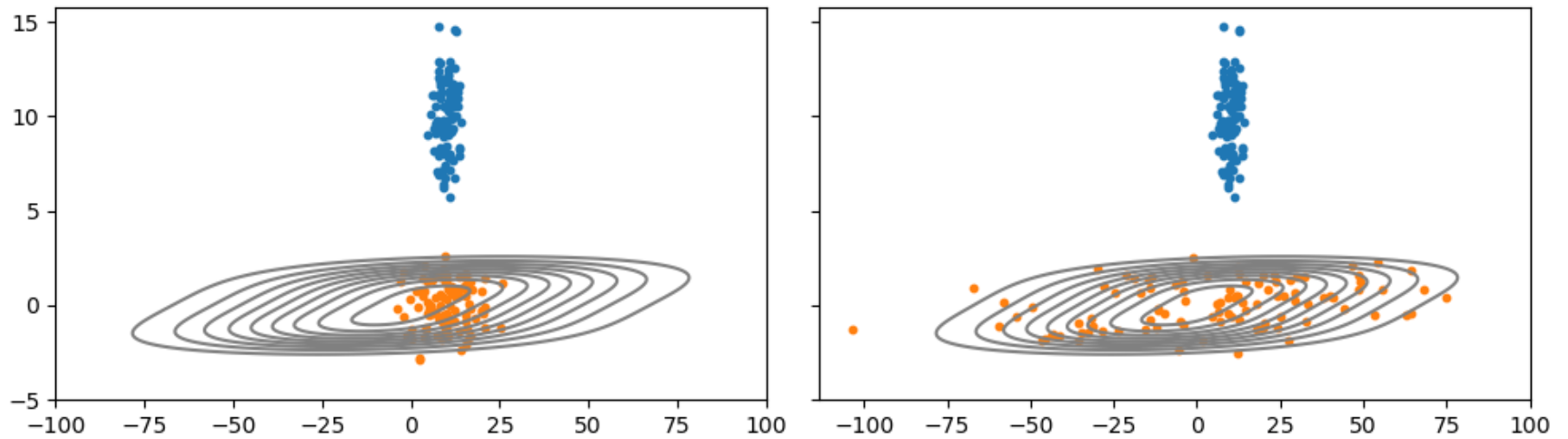
for preconditioner = 1:2
    theta = theta0
    for i = 1:N_t
        theta = update_ensemble(theta, dt, gradPhiR, preconditioner == 2)
    end

    ax[preconditioner].scatter(theta0[1, :], theta0[2, :], s = 10)
    ax[preconditioner].scatter(theta[1, :], theta[2, :], s = 10)

end

fig.tight_layout()
fig.savefig("GF.png")

```



In []: