

1. 指数分布族

对于指数分布族

$$\rho(\theta; a) = h(\theta) \exp\{T(\theta) \cdot a - A(a)\}$$
$$0 = \int \nabla_a \rho(\theta; a) da = \mathbb{E}_\rho[T(\theta) - \nabla_a A(a)]$$

Here we introduce another parameter vector λ . Secondly, the Fisher information matrix *FIM* is

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

2. 高斯变分推理

对于高斯分布

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

我们有

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

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

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

3. 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

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

它满足

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

因此

$$(XY)_{ij} = \sum_{l, m} X[i, j, l, m] Y[l, m]$$
$$= \frac{1}{4} \sum_{l, m} (C^{-1}[i, l] Y[l, m] C^{-1}[j, m] + C^{-1}[i, m] Y[l, m] C^{-1}[j, l])$$
$$= \frac{1}{4} (C^{-1} Y C^{-1} + C^{-1} Y^T C^{-1})_{ij}.$$

我们用了

由于

$$\frac{dC_i^{-1}}{dt} = -C_i^{-1} \frac{dC_i}{dt} C_i^{-1}$$

自然梯度下降法可以写成

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

对于协方差, 我们有

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

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

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

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, dt, N_t)
    N_0 = length(m_0)
    m = zeros(N_t+1, N_0)
    C = zeros(N_t+1, N_0, N_0)

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

    for i = 1:N_t
        EV0, EV0T = expectation(method_type, m_oo, C_oo, m[i, :], C[i, :, :])
        if method_type == "gradient_descent"
            m[i+1, :] = m[i, :] - EV0 * dt
            C[i+1, :, :] = C[i, :, :] + 1/2.0 * (inv(C[i, :, :]) - EV0T * dt)
        elseif method_type == "natural_gradient_descent"
            m[i+1, :] = m[i, :] - C[i, :, :] * EV0 * dt
            C[i+1, :, :] = inv(inv(C[i, :, :]) + EV0T - inv(C[i, :, :])) * dt)
        else
            error("Problem type : ", problem_type, " has not implemented!")
        end
    end
    return m, C
end

dt, 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, dt, N_t)
m_ngd, C_ngd = Continuous_Dynamics("natural_gradient_descent", m_oo, C_oo, m_0, C_0, dt, N_t)

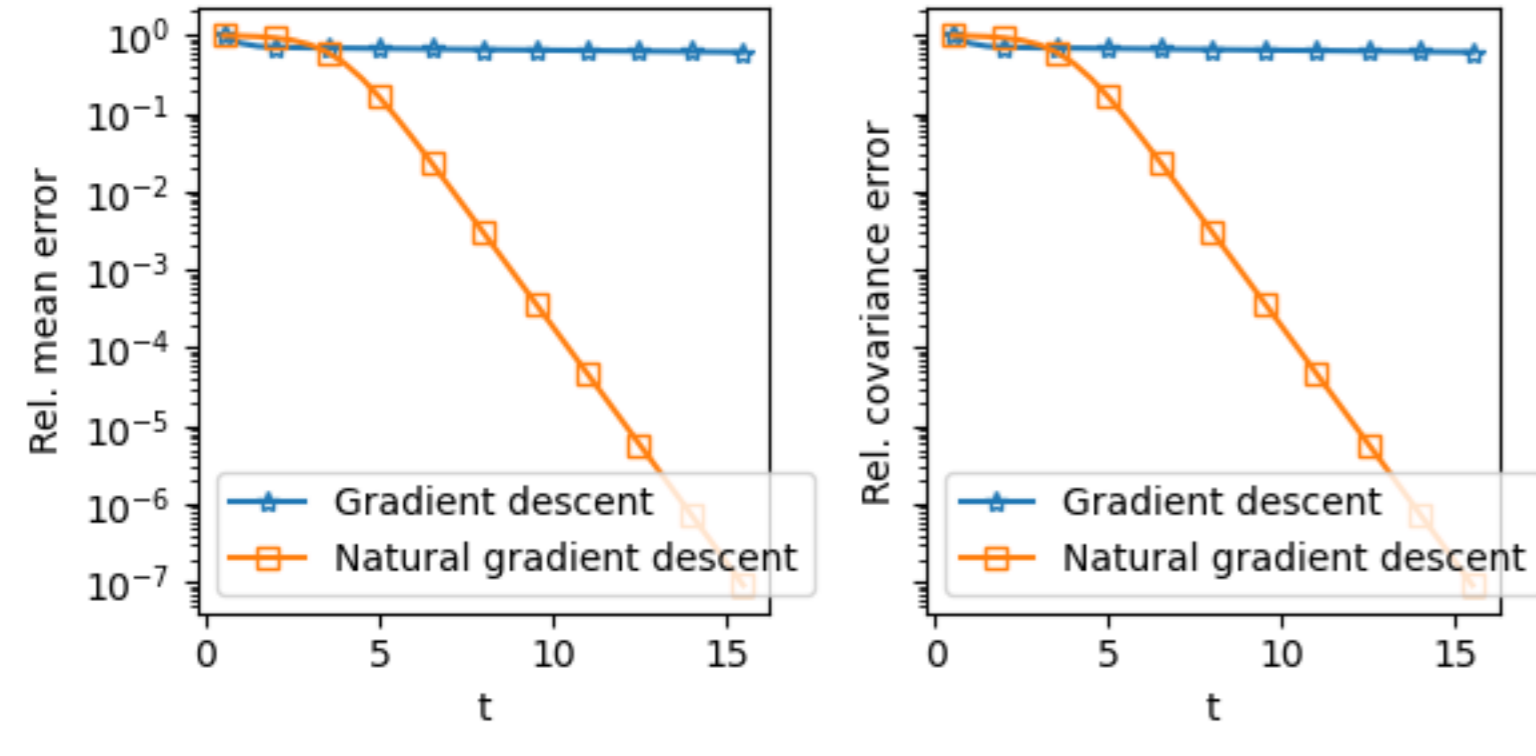
e_gd, e_ngd = zeros(N_t+1, 2), zeros(N_t+1, 2)

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} * dt
ax[1].semilogy(ts, e_gd[:, 1], "--", fillstyle="none", marker="x", label="Gradient descent")
ax[1].semilogy(ts, e_gd[:, 2], "-s", fillstyle="none", marker="x", label="Natural gradient descent")
ax[2].semilogy(ts, e_gd[:, 1], "--", fillstyle="none", marker="x", label="Gradient descent")
ax[2].semilogy(ts, e_gd[:, 2], "-s", fillstyle="none", marker="x", label="Natural gradient descent")

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")
```



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

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

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 $\tilde{\theta} = T\theta = A\theta + b$.

$$\tilde{\rho}(\tilde{\theta}) = \rho(\theta) |A^{-1}| \quad \tilde{\rho}_i(\tilde{\theta}) = \rho_i(\theta) |A^{-1}| \quad \tilde{C}_i = AC_i A^T$$
$$\nabla_{\tilde{\theta}} f(\tilde{\theta}) = A^{-T} \nabla_\theta f(\theta) \quad \nabla_{\tilde{\theta}} \cdot v(\tilde{\theta}) = \nabla_\theta \cdot (A^{-1} v(\theta))$$

带入Kalman-Wasserstein精度流

$$\frac{\partial \tilde{\rho}_i}{\partial t} = \nabla_{\tilde{\theta}} \cdot [\tilde{\rho}_i \tilde{C}_i (\nabla_{\tilde{\theta}} \log \tilde{\rho}^* - \nabla_{\tilde{\theta}} \log \tilde{\rho}_i)]$$
$$\Leftrightarrow \frac{\partial \rho_i}{\partial t} |A^{-1}| = \nabla_{\tilde{\theta}} \cdot [\rho_i AC_i A^T (\nabla_{\tilde{\theta}} \log \tilde{\rho}^* - \nabla_{\tilde{\theta}} \log \tilde{\rho}_i)] |A^{-1}|$$
$$\Leftrightarrow \frac{\partial \rho_i}{\partial t} |A^{-1}| = \nabla_\theta \cdot [\rho_i C_i (\nabla_\theta \log \tilde{\rho}^* - \nabla_\theta \log \tilde{\rho}_i)] |A^{-1}|$$
$$\Leftrightarrow \frac{\partial \rho_i}{\partial t} = \nabla_\theta \cdot [\rho_i C_i (\nabla_\theta \log \rho^* - \nabla_\theta \log \rho_i)]$$

6. 练习

In [2]:

```
using PyPlot, Random, Statistics, Distributions, ForwardDiff, LinearAlgebra

Random.seed!(42)

function Phi_R_Convex(theta)
    theta = 0.1
    theta_1, theta_2 = 0
    return (exp(theta_1 - theta_2)^2 + theta_2^4) / 20
end

Phi_R = Phi_R_Convex(theta)
function Phi_R(theta)
    return Phi_R(theta), ForwardDiff.gradient(Phi_R, theta)
end

function update_ensemble(theta, dt, func_neg_log_rho_i::Function, preconditioner::Bool)
    N_0, N_ens = size(theta)
    Phi_R, Phi_R_grad = zeros(N_ens), zeros(N_0, N_ens)
    for i = 1:N_ens
        Phi_R[i, :] = func_neg_log_rho_i(theta[i, :])
    end
    theta_mean = mean(theta, dims=2)
    theta_cov = ((theta - theta_mean) * (theta - theta_mean)^T) / N_ens
    Prec = (preconditioner ? theta_cov : I)
    O = (preconditioner ? chol(cholesky(Hermitian(theta_cov)).L : I)
    noise = rand(Normal{0, 1}, (N_0, N_ens))
    theta = theta - dt * Prec * Phi_R_grad + sqrt(2*dt) * (O * noise)
    return theta
end

update_ensemble (generic function with 1 method)
```

In [3]:

```
N_0 = 2
theta = [10.0; 10.0]
theta_0 = 0
theta_1 = [0.0; 0.0; 0.0 0.0]

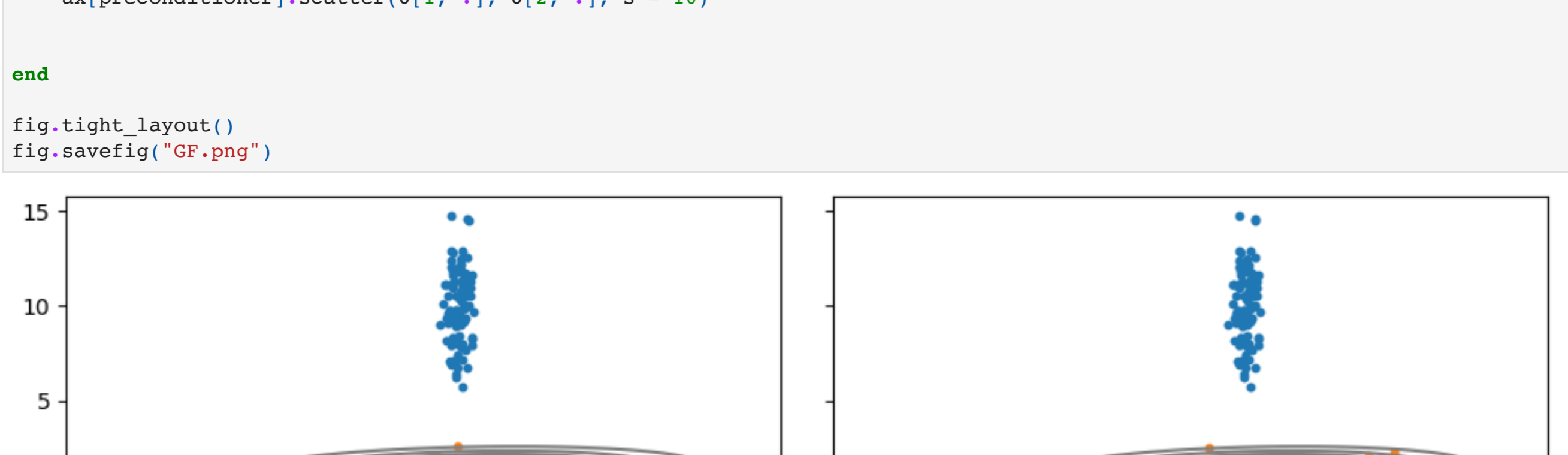
Lx, Ux = -100.0, 100.0
Ly, Uy = -5.0, 5.0
N = 500
X = zeros(N, N)
Y = zeros(N, N)
p = zeros(N, N)
for j = 1:N
    for i = 1:N
        X[i, j], Y[i, j] = Li * (Xj - Lx) + (i-1)/(N-1), Ly + (Uy - Ly) * (j-1)/(N-1)
    end
end
Z = sum(p)
p .= exp.(-p) / Z

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

N_ens = 100
theta = Array{rand(MvNormal{m_0, C_0}, N_ens)}
dt = 0.01
N_t = 2000

for preconditioner = 1:2
    theta = theta
    for i = 1:N_t
        theta = update_ensemble(theta, dt, Phi_R_grad, preconditioner == 2)
    end
    ax[preconditioner].scatter(theta[1, :], theta[2, :], s = 10)
    ax[preconditioner].scatter(theta[1, :], theta[2, :], s = 10)
end

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



In [1]: