

Langevin Dynamics for sampling and global optimization

Kirill Neklyudov





Goals of this talk

- Introduction to the Langevin dynamics
- Derive basics (1st half)
- Outline some important results (2nd half)
- Recommend some literature

Langevin Equation

Ito Stochastic Differential Equation (SDE):

$$\boxed{dX(t)} = -\nabla U(X(t))dt + \sigma \boxed{dBt}$$

 Force  Random fluctuations

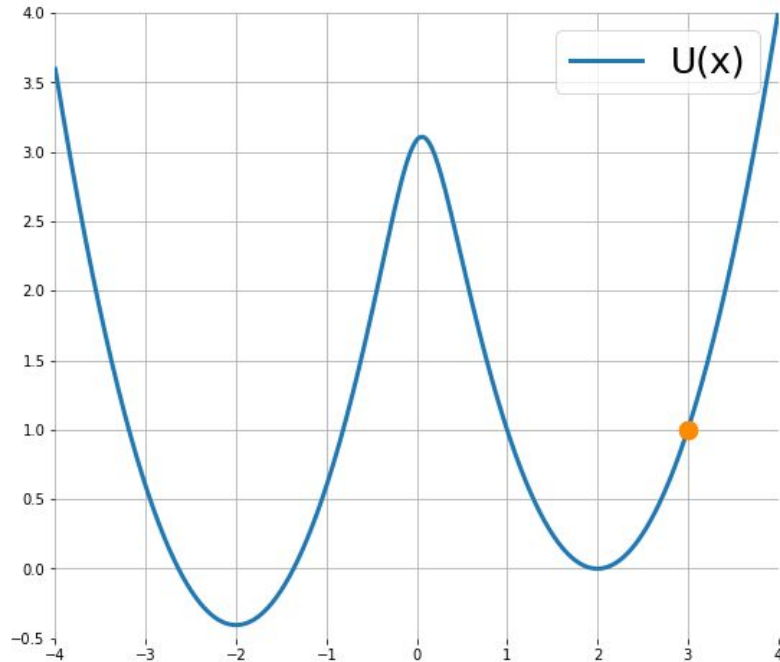
Discrete approximation:

$$\boxed{X_{t+1} - X_t} = -dt \nabla U(X_t) + \sigma \boxed{\sqrt{dt} \mathcal{N}(0, 1)}$$

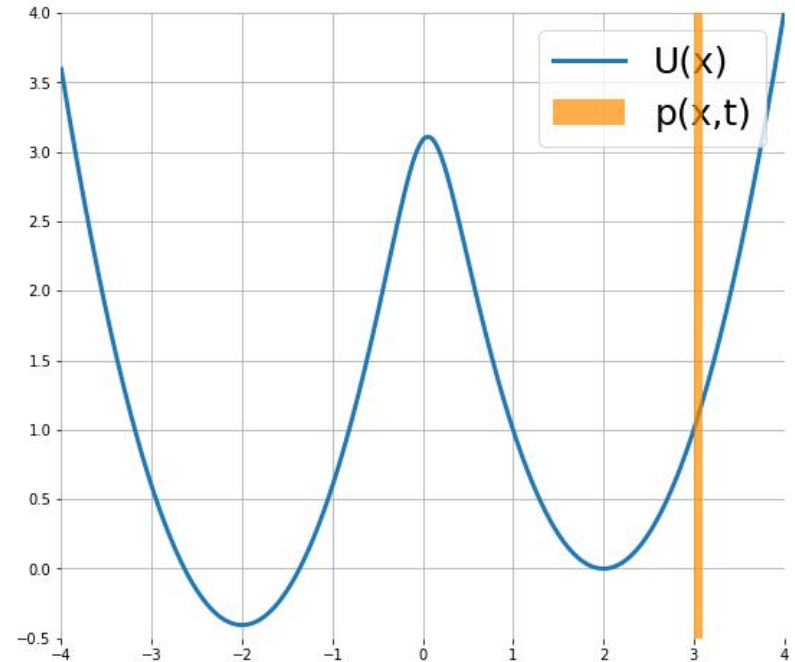
$$W_{t+1} - W_t = -\varepsilon \nabla \mathcal{L}(W_t) + \sigma \sqrt{\varepsilon} \mathcal{N}(0, 1)$$

1-d simulation

Langevin equation



Fokker-Planck equation



Derivation of the Fokker-Planck equation

Langevin equation:

$$dX(t) = -\nabla U(X(t))dt + \sigma dB_t$$

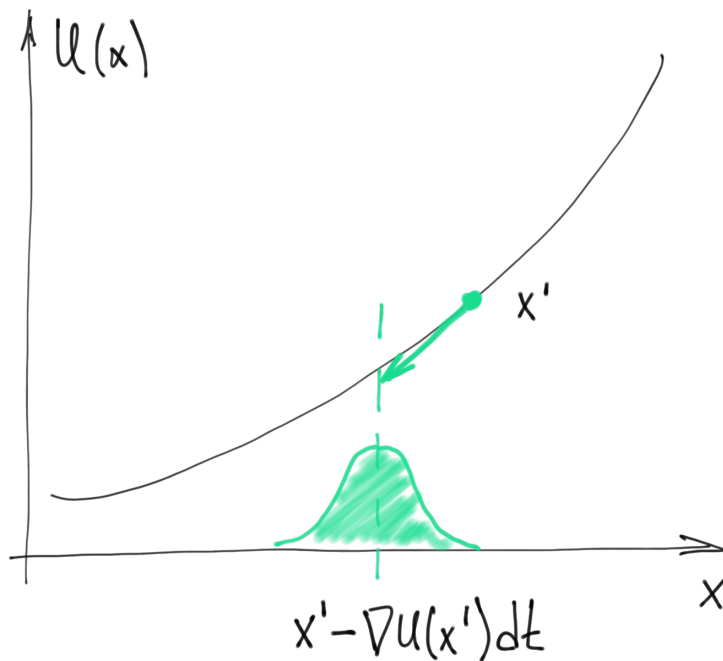
Increments of the Brownian motion:

$$dB_t \sim \mathcal{N}(0, dt \cdot I)$$

Consider a small increment of $X(t)$:

$$x - x' = -dt \nabla U(x') + \mathcal{N}(0, \sigma^2 dt)$$

$$x \sim \mathcal{N}\left(x' - \nabla U(x')dt, \sigma^2 dt\right)$$



Derivation of the Fokker-Planck equation

Density of particle distribution:

Unknown density!



$$p(x, t) = \int dx' p(x, t | x', t - dt) p(x', t - dt)$$

$$p(x, t | x', t - dt) = \frac{1}{(2\pi\sigma^2 dt)^{n/2}} \exp \left(\frac{-\overbrace{(x' - x - \nabla U(x')dt)^2}^y}{2\sigma^2 dt} \right)$$

Using the change of variables formula, we obtain:

$$p(x, t) = \int dy \left| \frac{\partial x'}{\partial y} \right| \mathcal{N}(y | 0, \sigma^2 dt \cdot I) p(x'(y), t - dt)$$

Changing variables

The change of variables:

$$y = x' - x - \nabla U(x')dt \quad \left| \frac{\partial x'}{\partial y} \right| = ? \quad x'(y) = ?$$



We can't invert this

$$y = x' - x - \left(\nabla U(x) + \frac{\partial \nabla U(x)}{\partial x} (x' - x)dt + o(x' - x) \right) dt$$

$$\left(I - \frac{\partial \nabla U(x)}{\partial x} dt \right) x' = y + x + \nabla U(x)dt - \frac{\partial \nabla U(x)}{\partial x} xdt + o(dt)$$

$$x' = \left(I - \frac{\partial \nabla U(x)}{\partial x} dt \right)^{-1} \left(y + x + \nabla U(x)dt - \frac{\partial \nabla U(x)}{\partial x} xdt + o(dt) \right)$$

Changing variables

$$\begin{aligned}
 x' &= \left(I - \frac{\partial \nabla U(x)}{\partial x} dt \right)^{-1} \left(y + x + \nabla U(x) dt - \frac{\partial \nabla U(x)}{\partial x} x dt + o(dt) \right) \\
 &= \left(I + \frac{\partial \nabla U(x)}{\partial x} dt + o(dt) \right) \left(y + x + \nabla U(x) dt - \frac{\partial \nabla U(x)}{\partial x} x dt + o(dt) \right) \\
 &= y + x + \nabla U(x) dt - \frac{\partial \nabla U(x)}{\partial x} x dt + \frac{\partial \nabla U(x)}{\partial x} y dt + \frac{\partial \nabla U(x)}{\partial x} x dt + o(dt) \\
 &= x + y + \nabla U(x) dt + \frac{\partial \nabla U(x)}{\partial x} y dt + o(dt)
 \end{aligned}$$

Changing variables

From the previous slide:

$$x' = x + y + \nabla U(x)dt + \frac{\partial \nabla U(x)}{\partial x} \boxed{ydt} + o(dt)$$

$$ydt = (x' - x - \nabla U(x)dt)dt$$

$$x - x' = -dt \nabla U(x') + \mathcal{N}(0, \sigma^2 dt)$$

Equation for the increment
(Langevin equation)

$$ydt = dt \underbrace{\sqrt{dt} \mathcal{N}(0, \sigma^2)}_{dB_t} = o(dt)$$

Note that

dB_t

$$y = \sqrt{dt} \mathcal{N}(0, \sigma^2) \neq o(dt) \quad \text{since} \quad \lim_{dt \rightarrow 0} \frac{\sqrt{dt} \mathcal{N}(0, \sigma^2)}{dt} = \infty$$

Changing variables

Finally!

$$x' = x + y + \nabla U(x)dt + o(dt)$$

With a little more efforts (homework):

$$\left| \frac{\partial x'}{\partial y} \right| = 1 + \operatorname{div} \nabla U(x)dt + o(dt)$$

Reminder: $\operatorname{div} \vec{f}(\vec{x}) = \frac{\partial f_1}{\partial x_1} + \frac{\partial f_2}{\partial x_2} + \dots + \frac{\partial f_n}{\partial x_n}$

Derivation of the Fokker-Planck equation

$$p(x, t) = \int dy \left| \frac{\partial x'}{\partial y} \right| \mathcal{N}(y|0, \sigma^2 dt \cdot I) p(x'(y), t - dt) \quad \text{Increment of the density}$$

$$\left. \begin{aligned} x' &= x + y + \nabla U(x) dt + o(dt) \\ \left| \frac{\partial x'}{\partial y} \right| &= 1 + \operatorname{div} \nabla U(x) dt + o(dt) \end{aligned} \right\} \quad \text{Change of variables}$$

$$p(x, t) = (1 + \operatorname{div} \nabla U(x) dt) \mathbb{E}_y \left[p(y + x + \nabla U(x) dt, t - dt) \right] \quad \text{???}$$
$$y \sim \mathcal{N}(0, \sigma^2 dt \cdot I)$$

Derivation of the Fokker-Planck equation

$$\mathbb{E}_y \left[p(y + x + \nabla U(x)dt, t - dt) \right] = \mathbb{E}_y \left[$$

0th order: $p(x, t) +$

1st order: $+ \nabla_x p(x, t)(y + \nabla U(x)dt) + \frac{\partial}{\partial t} p(x, t)(-dt) +$

2nd order: $+ \frac{1}{2}(y + \nabla U(x)dt)^T \frac{\partial^2 p(x, t)}{\partial x^2} (y + \nabla U(x)dt) \Big]$

Taking the expectation

0th order: $p(x, t) +$

$$\mathbb{E}_y p(x, t) = p(x, t)$$

1st order: $+ \nabla_x p(x, t)(y + \nabla U(x)dt) + \frac{\partial}{\partial t} p(x, t)(-dt) +$

$$\mathbb{E}_y \left[\nabla_x p(x, t)^T (y + \nabla U(x)dt) \right] = \nabla_x p(x, t)^T \mathbb{E}_y[y] + dt \nabla_x p(x, t)^T \nabla U(x)$$

$$y \sim \mathcal{N}(0, \sigma^2 dt)$$

$$= 0 + dt \nabla_x p(x, t)^T \nabla U(x)$$

$$\mathbb{E}_y \left[dt \frac{\partial}{\partial t} p(x, t) \right] = dt \frac{\partial}{\partial t} p(x, t)$$

Taking the expectation (2nd order)

$$\begin{aligned}
 & \mathbb{E}_y \left[(y + \nabla U(x)dt)^T \frac{\partial^2 p(x, t)}{\partial x^2} (y + \nabla U(x)dt) \right] = \\
 &= \mathbb{E}_y \left[y^T \frac{\partial^2 p(x, t)}{\partial x^2} y + 2dt \nabla U(x)^T \frac{\partial^2 p(x, t)}{\partial x^2} y + dt^2 \nabla U(x)^T \frac{\partial^2 p(x, t)}{\partial x^2} \nabla U(x) \right] = \\
 &= \mathbb{E}_y \left[\sum_{i,j} \left(\frac{\partial^2 p(x, t)}{\partial x^2} \right)_{ij} y_i y_j \right] + 2dt \nabla U(x)^T \frac{\partial^2 p(x, t)}{\partial x^2} \mathbb{E}_y y + o(dt) = \\
 &= \sum_{i=j} \left(\frac{\partial^2 p(x, t)}{\partial x^2} \right)_{ii} \mathbb{E}_y [y_i^2] + \sum_{i \neq j} \left(\frac{\partial^2 p(x, t)}{\partial x^2} \right)_{ij} \mathbb{E}_y [y_i y_j] + o(dt) = \\
 &\quad = \Delta p(x, t) = \sigma^2 dt \qquad \qquad \qquad = 0
 \end{aligned}$$

Derivation of the Fokker-Planck equation

$$p(x, t) = (1 + \operatorname{div} \nabla U(x) dt) \mathbb{E}_y \left[p(y + x + \nabla U(x) dt, t - dt) \right] \quad \text{Increment of the density}$$

$$p(x, t) = (1 + \operatorname{div} \nabla U(x) dt) \left(p(x, t) + dt \nabla_x p(x, t)^T \nabla U(x) - dt \frac{\partial}{\partial t} p(x, t) + \frac{1}{2} \sigma^2 dt \Delta p(x, t) + o(dt) \right) \quad \text{Taylor series}$$

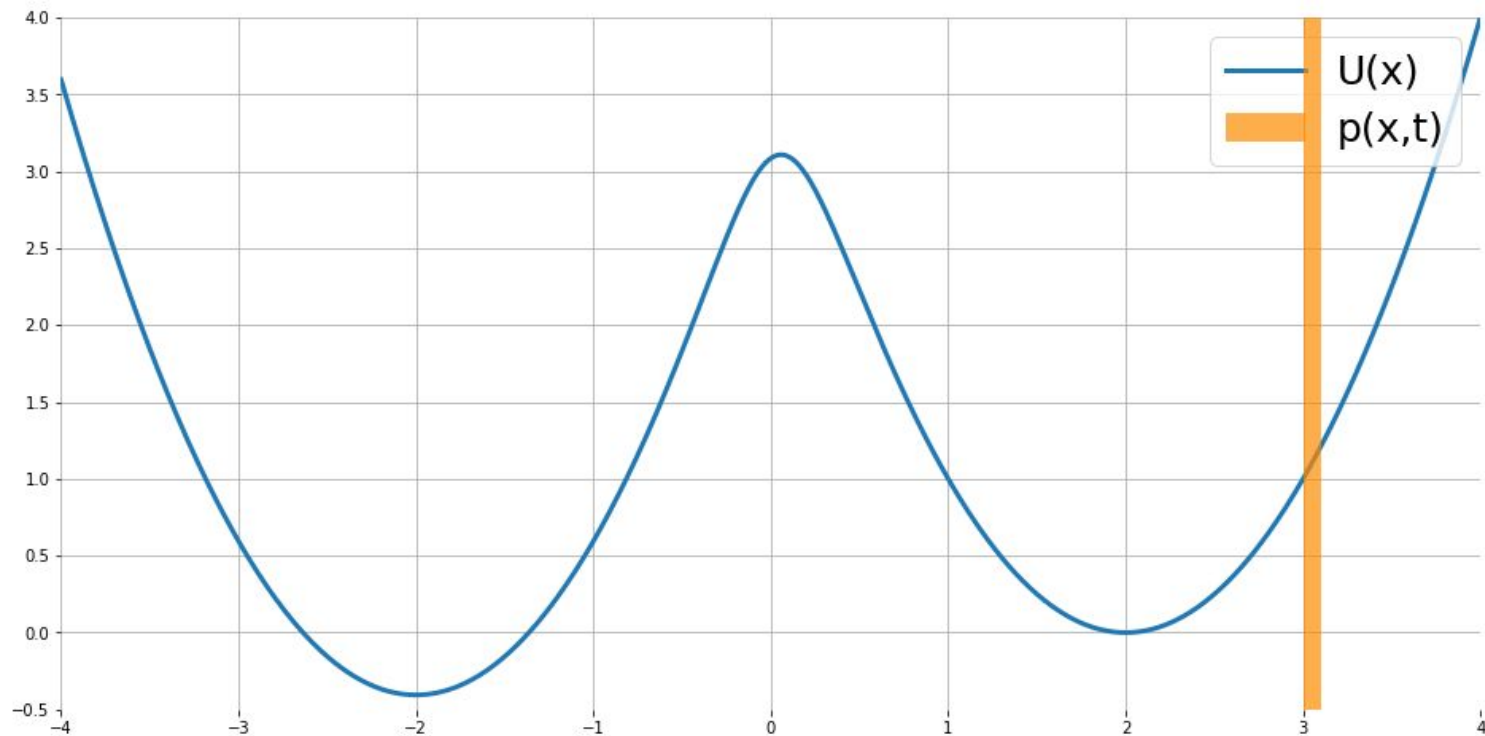
$$\cancel{p(x, t)} = \cancel{p(x, t)} + dt \nabla_x p(x, t)^T \nabla U(x) - dt \frac{\partial}{\partial t} p(x, t) + \frac{1}{2} \sigma^2 dt \Delta p(x, t) + p(x, t) \operatorname{div} \nabla U(x) dt + o(dt)$$

$$\frac{\partial}{\partial t} p(x, t) = \nabla_x p(x, t)^T \nabla U(x) + p(x, t) \operatorname{div} \nabla U(x) + \frac{1}{2} \sigma^2 \Delta p(x, t) + \cancel{\frac{o(dt)}{dt}} \quad 0$$

5 min break

Fokker-Planck equation

$$\frac{\partial}{\partial t}p(x,t) = \nabla_x p(x,t)^T \nabla U(x) + p(x,t) \operatorname{div} \nabla U(x) + \frac{1}{2} \sigma^2 \Delta p(x,t)$$



Stationary distribution of the Langevin dynamics

$$\frac{\partial}{\partial t} p(x, t) = \nabla_x p(x, t)^T \nabla U(x) + p(x, t) \operatorname{div} \nabla U(x) + \frac{1}{2} \sigma^2 \Delta p(x, t)$$

$$p(x, t) = \hat{p}(x) \quad (\text{density does not change anymore})$$

$$\text{Let } p_G(x) = \frac{1}{Z} \exp \left(-\frac{U(x)}{T} \right), \quad Z = \int dx \exp \left(-\frac{U(x)}{T} \right) \quad \text{Gibbs distribution}$$

With a little efforts (homework), we obtain:

$$0 = \nabla \hat{p}(x)^T \nabla U(x) + \hat{p}(x) \operatorname{div} \nabla U(x) + \frac{1}{2} \sigma^2 \Delta \hat{p}(x) \quad \text{when} \quad T = \frac{\sigma^2}{2}$$

Sampling via the Langevin dynamics

$$dX(t) = -\nabla U(X(t))dt + \sigma dBt \quad \text{Langevin equation}$$

Particles have the stationary distribution:

$$p_G(x) = \frac{1}{Z} \exp\left(-\frac{2U(x)}{\sigma^2}\right), \quad Z = \int dx \exp\left(-\frac{2U(x)}{\sigma^2}\right) \quad \text{Gibbs distribution}$$

We want to sample from $p(x) = \frac{\hat{p}(x)}{Z'}$

$$U(x) = -\log p(x), \quad \sigma = \sqrt{2}$$

$$p_G(x) = \frac{1}{Z} \exp\left(-\frac{-2\log p(x)}{2}\right) = p(x), \quad Z = \int dx \exp(\log p(x)) = 1$$

Note! $\nabla U(x) = -\nabla \log p(x) = -\nabla \log \hat{p}(x) - \nabla \log \cancel{Z'}$

Sampling via the Langevin dynamics

Stochastic Differential Equation for sampling:

$$dX(t) = \nabla \log p(X(t))dt + \sqrt{2}dB_t$$

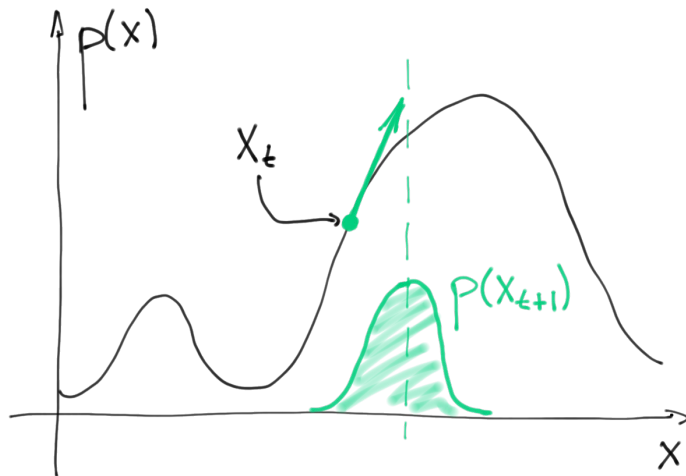
Discrete approximation:

$$X_{t+1} = X_t + dt \nabla \log p(X_t) + \mathcal{N}(0, 2dt)$$

More popular way:

$$X_{t+1} = X_t + \frac{\varepsilon}{2} \nabla \log p(X_t) + \mathcal{N}(0, \varepsilon)$$

$$X_t \sim p(x), \quad \forall t > t_\infty$$



Langevin dynamics for the Bayesian inference

Predictive distribution

$$p(y|D_{\text{train}}) = \mathbb{E}_{p(\theta|D_{\text{train}})} p(y|\theta) \simeq \frac{1}{K} \sum_{i=1}^K p(y|\theta_i), \quad \theta_i \sim p(\theta|D_{\text{train}})$$

We need samples

$$d\theta(t) = \nabla \log p(\theta(t)|D_{\text{train}}) dt + \sqrt{2} dB_t$$

Langevin equation

???

$$\theta_{t+1} = \theta_t + \frac{\varepsilon}{2} \nabla \log p(\theta_t|D_{\text{train}}) + \mathcal{N}(0, \varepsilon)$$

Discrete approximation

$$\theta_{t+1} = \theta_t + \frac{\varepsilon}{2} \nabla_{\theta} \left(\sum_{i=1}^N \log p(\theta_t|(x_i, y_i)) + \log p(\theta_t) \right) + \mathcal{N}(0, \varepsilon)$$

???

$$\theta_{t+1} = \theta_t + \frac{\varepsilon}{2} \nabla_{\theta} \left(\frac{N}{B} \sum_{k=1}^B \log p(\theta_t|(x_{i_k}, y_{i_k})) + \log p(\theta_t) \right) + \mathcal{N}(0, \varepsilon)$$

Borkar, Mitter, 1999

Consider a SDE:

$$dX(t) = h(X(t))dt + \sigma dBt$$

Stationary distribution

$$X(t) \sim p_\sigma(x), \quad t > t_\infty$$

Discrete approximation:

$$X_{k+1} = X_k + \varepsilon(h(X_k) + M_k) + \mathcal{N}(0, \sigma^2 \varepsilon) \quad X_k \sim \hat{p}_\sigma(x), \quad k > k_\infty$$

Stationary distribution

$$\mathbb{E}M_k = 0, \quad \forall k$$

Theorem

$$\forall \delta > 0, \exists \varepsilon : \text{KL}(p_\sigma(x) || \hat{p}_\sigma(x)) < \delta$$

Sketch of the proof

$$\tilde{X}(t) = X(0) + \int_0^t \left(h(\tilde{X}(\lfloor s \rfloor_\varepsilon) + \xi_s) \right) ds + \sigma \tilde{B}(t)$$

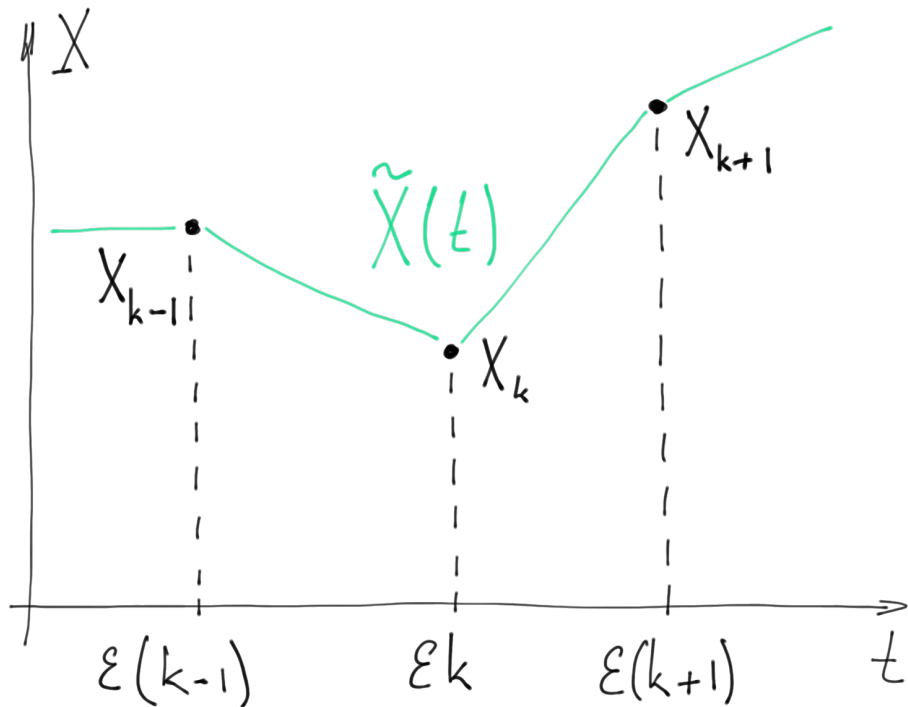
$$\lfloor s \rfloor_\varepsilon = k\varepsilon, \text{ if } s \in [k\varepsilon, (k+1)\varepsilon)$$

$$\xi_s = M_k, \text{ if } s \in [k\varepsilon, (k+1)\varepsilon)$$

$$\tilde{B}((k+1)\varepsilon) - \tilde{B}(k\varepsilon) = \mathcal{N}(0, \varepsilon)$$

Lemma

$$\forall t \quad \mathbb{E} \left[\|X(t) - \tilde{X}(t)\|^2 \right] \rightarrow 0, \text{ as } \varepsilon \rightarrow 0$$



What happened to the noise?

$$dX(t) = h(X(t))dt + \sigma dBt \quad \text{Original dynamics}$$

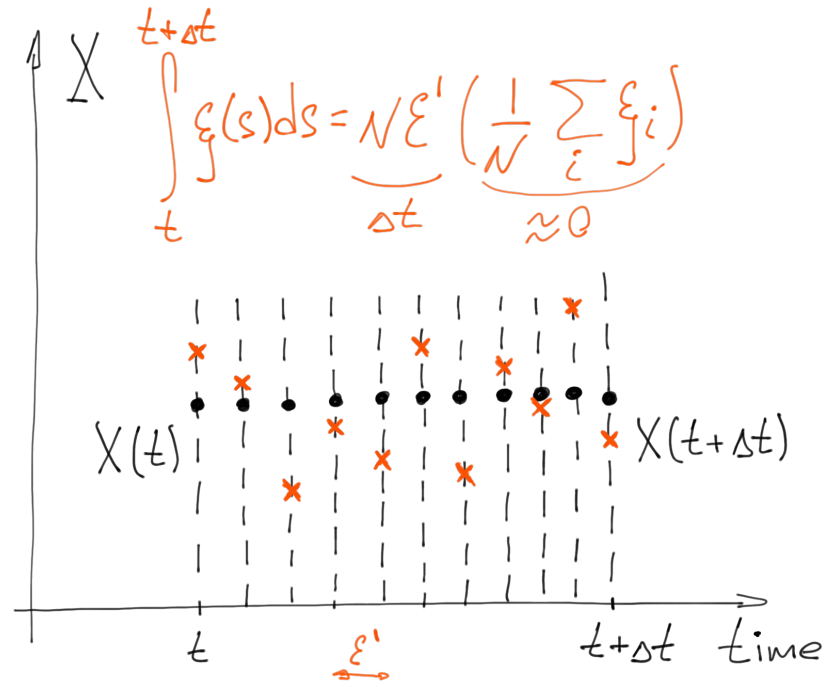
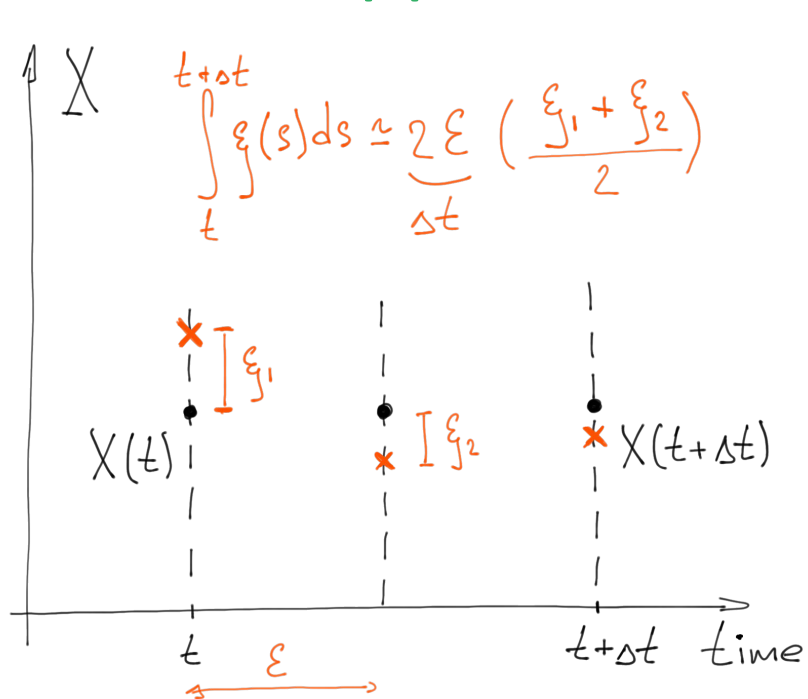
$$X(t) = X(0) + \int_0^t h(X(s))ds + \sigma B(t) \quad \text{The same but integrated}$$

$$\tilde{X}(t) = X(0) + \int_0^t \left(h(\tilde{X}(\lfloor s \rfloor_\varepsilon) + \xi_s) \right) ds + \sigma \tilde{B}(t) \quad \text{Our approximation}$$

Free speed-up?

Lemma says $X(t) = \tilde{X}(t)$, when $\varepsilon \rightarrow 0$

What happened to the noise?



Computational efforts are hidden here

$$X(t) = \tilde{X}(t), \text{ when } \boxed{\varepsilon \rightarrow 0}$$

Sketch of the proof

Stationary distribution

$$X(t) \sim p_\sigma(x), \quad t > t_\infty$$

Stationary distribution

$$X_k \sim \hat{p}_\sigma(x), \quad k > k_\infty$$

Lemma 1

$$\forall t \quad \mathbb{E} \left[\|X(t) - \tilde{X}(t)\|^2 \right] \rightarrow 0, \quad \text{as } \varepsilon \rightarrow 0$$

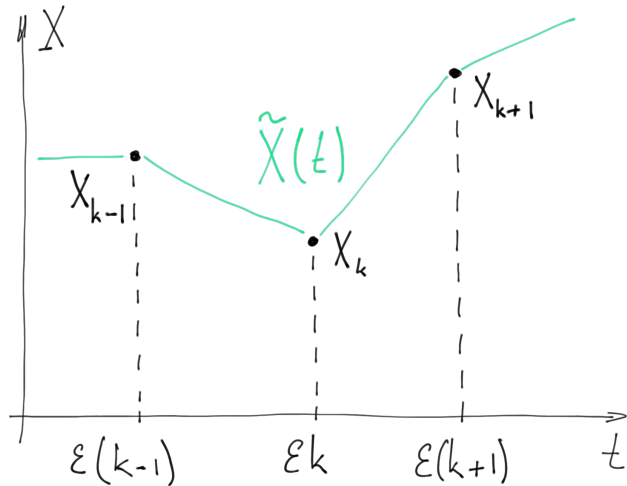
Lemma 2

$$X(t) \sim p(x, t)$$

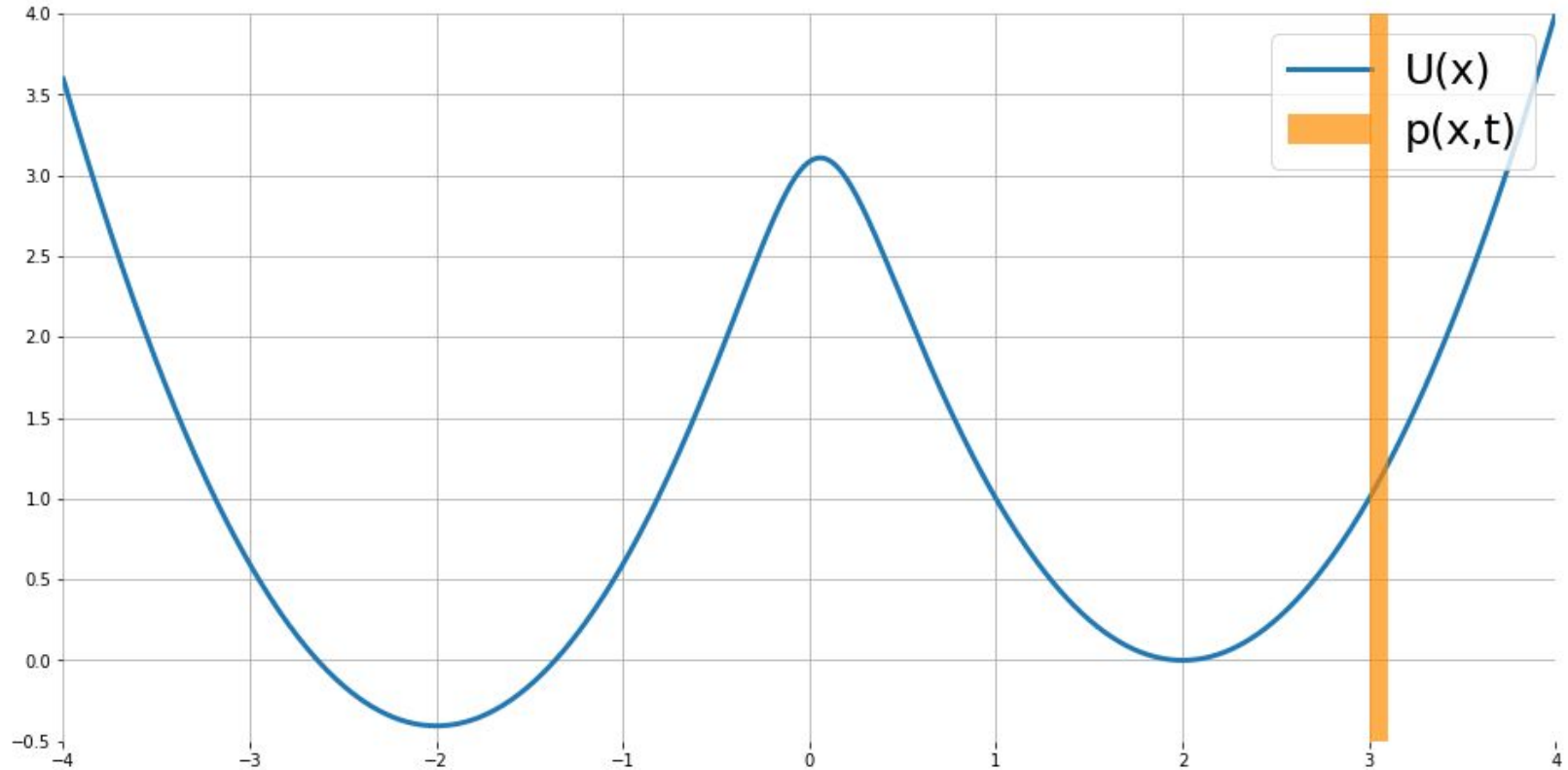
$\text{KL}(p_\sigma(x) || p(x, t))$ is strictly decreasing in t

Theorem

$$\forall \delta > 0, \exists \varepsilon : \text{KL}(p_\sigma(x) || \hat{p}_\sigma(x)) < \delta$$



Global optimization



Temperature annealing

$$dX(t) = -\nabla U(X(t))dt + \sigma dB_t \quad \text{Langevin equation}$$

Particles have the stationary distribution:

$$p_T(x) = \frac{1}{Z} \exp\left(-\frac{U(x)}{T}\right), \quad Z = \int dx \exp\left(-\frac{U(x)}{T}\right) \quad \text{Gibbs distribution} \quad T = \frac{\sigma^2}{2}$$

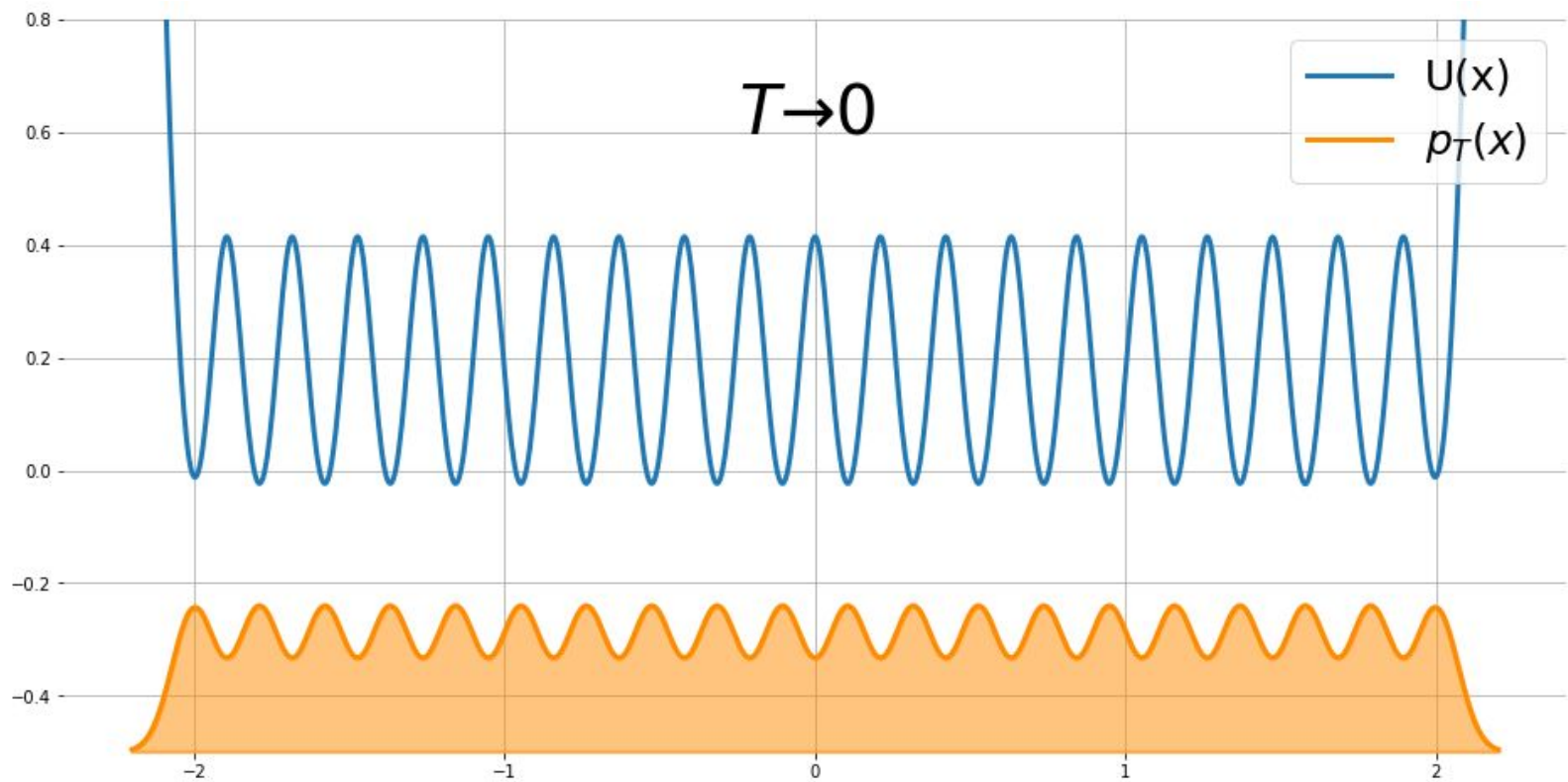
$$dX(t) = -\nabla U(X(t))dt + \sqrt{2T}dB_t$$

$$p_T(x) = \exp\left(-\frac{U(x)}{T}\right) \xrightarrow[T \rightarrow 0]{\mathcal{D}} \pi(x)$$



Concentrates on the global minima

Annealing example



DIFFUSION FOR GLOBAL OPTIMIZATION IN \mathbb{R}^n *

TZUU-SHUH CHIANG[†], CHII-RUEY HWANG[†] AND SHUENN-JYI SHEU[†]

$$dX(t) = -\nabla U(X(t))dt + \sqrt{2T(t)}dB_t \quad \text{Langevin equation}$$

???

$$X(t) \sim p(x, t) \quad \text{Distribution of particles}$$

$$p_T(x) = \exp\left(-\frac{U(x)}{T}\right) \xrightarrow{T \rightarrow 0} \pi(x)$$

Theorem:

Annealing schedule

$$T(t) = \frac{c}{\log t}$$

$$p(x, t) \xrightarrow[t \rightarrow \infty]{P} \pi(x)$$

RECURSIVE STOCHASTIC ALGORITHMS FOR GLOBAL OPTIMIZATION IN \mathbb{R}^{d*}

SAUL B. GELFAND[†] AND SANJOY K. MITTER[‡]

$$X_{k+1} = X_k - \varepsilon_k(\nabla U(X_k) + M_k) + \sqrt{2T_k}\mathcal{N}(0, 1) \quad \mathbb{E}M_k = 0, \quad \forall k$$

$$X_k \sim p(x, k) \quad \text{Distribution of particles} \quad p_T(x) = \exp\left(-\frac{U(x)}{T}\right) \xrightarrow[T \rightarrow 0]{\mathcal{D}} \pi(x)$$

Theorem:

Annealing schedule

$$\varepsilon_k = \frac{c_1}{k} \quad T_k = \frac{c_2}{k \log \log k} \quad p(x, k) \xrightarrow[t \rightarrow \infty]{P} \pi(x)$$

Non-Convex Learning via Stochastic Gradient Langevin Dynamics: A Nonasymptotic Analysis

Maxim Raginsky*

Alexander Rakhlin[†]

Matus Telgarsky[‡]

Continuous process

$$X(t) \sim p(x, t)$$

Discrete approximation

$$X_k \sim p(x, k)$$

Stationary distribution

$$p_T(x) = \frac{1}{Z} \exp \left(-\frac{U(x)}{T} \right)$$

Discretization error

Convergence speed

$$W_2 \left(p(x, k), p_T(x) \right) = W_2 \left(p(x, k), p(x, t) \right) + W_2 \left(p(x, t), p_T(x) \right)$$

$$W_2 \left(p(x, k), p_T(x) \right) = O \left((C + \varepsilon^{1/4}) \varepsilon k \right) + O \left(\exp(-\varepsilon k) \right)$$

Further reading

C.W. Gardiner

Handbook of Stochastic Methods

for Physics, Chemistry and the Natural Sciences

Second Edition
With 29 Figures



Springer

Bernt Øksendal

Stochastic Differential Equations

An Introduction with Applications
Fifth Edition, Corrected Printing
Springer-Verlag Heidelberg New York

Springer-Verlag
Berlin Heidelberg New York
London Paris Tokyo
Hong Kong Barcelona
Budapest

The end