# Advanced Computational Methods in Statistics
# Lecture 2

O. Deniz Akyildiz

LTCC Advanced Course

October 14, 2024

**IMPERIAL**

https://akyildiz.me/

𝕏: @odakyildiz

Recall our basic task:

Recall our basic task:

- We want to sample from a distribution $\pi(x) \propto \gamma(x)$ given only the knowledge of $\gamma(x)$.

Recall our basic task:

▶ We want to sample from a distribution $\pi(x) \propto \gamma(x)$ given only the knowledge of $\gamma(x)$.

▶ We want to use these samples to estimate an integral

$$(\varphi, \pi) = \int \varphi(x)\pi(x)\,\mathrm{d}x$$

Last week, we have covered the basic sampling techniques:

Last week, we have covered the basic sampling techniques:

- ▶ Uniform random number generation
  - ▶ Linear congruential generators

Last week, we have covered the basic sampling techniques:

- ▶ Uniform random number generation
  - ▶ Linear congruential generators
- ▶ Inversion (inverse transform) sampling
  - ▶ $U \sim \mathcal{U}(0,1)$
  - ▶ $X = F^{-1}(U)$

Last week, we have covered the basic sampling techniques:

- ▶ Uniform random number generation
  - ▶ Linear congruential generators
- ▶ Inversion (inverse transform) sampling
  - ▶ $U \sim \mathcal{U}(0, 1)$
  - ▶ $X = F^{-1}(U)$
- ▶ Rejection sampling
  - ▶ $X' \sim q(x)$
  - ▶ Accept $X'$ with probability $\gamma(X')/Mq(X')$

Last week, we have covered the basic sampling techniques:

- ▶ Uniform random number generation
    - ▶ Linear congruential generators
- ▶ Inversion (inverse transform) sampling
    - ▶ $U \sim \mathcal{U}(0,1)$
    - ▶ $X = F^{-1}(U)$
- ▶ Rejection sampling
    - ▶ $X' \sim q(x)$
    - ▶ Accept $X'$ with probability $\gamma(X')/Mq(X')$
- ▶ Importance sampling
    - ▶ Sample $X_1, \ldots, X_N \sim q(x)$
    - ▶ Estimate $(\varphi, \pi) \approx \sum_{i=1}^N \varphi(X_i) \mathrm{w}_i,$

Last week, we have covered the basic sampling techniques:

▶ Uniform random number generation
   ▶ Linear congruential generators
▶ Inversion (inverse transform) sampling
   ▶ $U \sim \mathcal{U}(0,1)$
   ▶ $X = F^{-1}(U)$
▶ Rejection sampling
   ▶ $X' \sim q(x)$
   ▶ Accept $X'$ with probability $\gamma(X')/Mq(X')$
▶ Importance sampling
   ▶ Sample $X_1, \ldots, X_N \sim q(x)$
   ▶ Estimate $(\varphi, \pi) \approx \sum_{i=1}^{N} \varphi(X_i) \mathrm{w}_i$,

The code is also available for these parts:

```
https://akyildiz.me/advanced-computational-statistics
```

OK, so what is wrong with these methods?

Let us exemplify a few issues. Consider the following target distribution on $\mathbb{R}^d$:

$$\pi(x) = \frac{1}{\sigma_\pi^d (2\pi)^{d/2}} \exp\left(-\frac{1}{2\sigma_\pi^2} \|x\|^2\right)$$

and the following proposal distribution:

$$q(x) = \frac{1}{\sigma_q^d (2\pi)^{d/2}} \exp\left(-\frac{1}{2\sigma_q^2} \|x\|^2\right)$$

where $\sigma_q > \sigma_\pi$.

We know that the acceptance probability is

$$\alpha(x) = \frac{\pi(x)}{Mq(x)}.$$

Mini-quiz: How do we choose $M$?

We know that the acceptance probability is

$$\alpha(x) = \frac{\pi(x)}{Mq(x)}.$$

Mini-quiz: How do we choose $M$?

$$M = \sup_{x \in \mathbb{R}^d} \frac{\pi(x)}{q(x)}.$$

We know that the acceptance probability is

$$\alpha(x) = \frac{\pi(x)}{Mq(x)}.$$

Mini-quiz: How do we choose $M$?

$$M = \sup_{x \in \mathbb{R}^d} \frac{\pi(x)}{q(x)}.$$

Then, we can write

$$M = \sup_{x \in \mathbb{R}^d} \frac{\sigma_q}{\sigma_\pi} \exp\left(-\frac{1}{2\sigma_\pi^2}\|x\|^2 + \frac{1}{2\sigma_q^2}\|x\|^2\right)$$

$$= \frac{\sigma_q^d}{\sigma_\pi^d} \sup_{x \in \mathbb{R}^d} \exp\left(\frac{\sigma_\pi^2 - \sigma_q^2}{2\sigma_q^2 \sigma_\pi^2}\|x\|^2\right) = \frac{\sigma_q^d}{\sigma_\pi^d}.$$

Mini-quiz: Given $M$, what is the acceptance rate?

Mini-quiz: Given $M$, what is the acceptance rate?

$$\hat{a} = \frac{1}{M} = \frac{\sigma_\pi^d}{\sigma_q^d}.$$

This means that as $d \to \infty$, given $\sigma_q > \sigma_\pi$, $\hat{a} \to 0$.

The curse of dimensionality for rejection samplers.

In standard Monte Carlo methods course, you would hear things like

In standard Monte Carlo methods course, you would hear things like

► Monte Carlo estimators are independent of the dimension of the problem.

In standard Monte Carlo methods course, you would hear things like

- ▶ Monte Carlo estimators are independent of the dimension of the problem.
- ▶ Importance sampling estimators are also independent of the dimension of the problem.

In standard Monte Carlo methods course, you would hear things like

▶ Monte Carlo estimators are independent of the dimension of the problem.

▶ Importance sampling estimators are also independent of the dimension of the problem.

These are false statements.

Importance sampling estimators also suffer badly as $d \to \infty$ (Li et al., 2005).

This motivates us to move on to our next topic: Markov chain Monte Carlo methods.

▶ In both high-dimensional sampling and more generally generative modelling, techniques based on MCMC and similar ideas are the state-of-the-art.

This motivates us to move on to our next topic: Markov chain Monte Carlo methods.

- ▶ In both high-dimensional sampling and more generally generative modelling, techniques based on MCMC and similar ideas are the state-of-the-art.
- ▶ Of course, there are many other techniques that are used in practice, but MCMC is the most popular one.

This motivates us to move on to our next topic: Markov chain Monte Carlo methods.

- ▶ In both high-dimensional sampling and more generally generative modelling, techniques based on MCMC and similar ideas are the state-of-the-art.
- ▶ Of course, there are many other techniques that are used in practice, but MCMC is the most popular one.

Next up: Introducing Markov chains.

A discrete-time Markov chain is a sequence of random variables $X_0, X_1, \ldots$ such that:

A discrete-time Markov chain is a sequence of random variables $X_0, X_1, \ldots$ such that:

- $X_t$ depends only on $X_{t-1}$

A discrete-time Markov chain is a sequence of random variables $X_0, X_1, \ldots$ such that:

- $X_t$ depends only on $X_{t-1}$
- In other words, $X_t$ is conditionally independent of $X_0, \ldots, X_{t-2}$ given $X_{t-1}$.

A discrete-time Markov chain is a sequence of random variables $X_0, X_1, \ldots$ such that:

- $X_t$ depends only on $X_{t-1}$
- In other words, $X_t$ is conditionally independent of $X_0, \ldots, X_{t-2}$ given $X_{t-1}$.

The evolution of the chain is governed by:

A discrete-time Markov chain is a sequence of random variables $X_0, X_1, \ldots$ such that:

- $X_t$ depends only on $X_{t-1}$
- In other words, $X_t$ is conditionally independent of $X_0, \ldots, X_{t-2}$ given $X_{t-1}$.

The evolution of the chain is governed by:

- A transition matrix $M$ (discrete case)

A discrete-time Markov chain is a sequence of random variables $X_0, X_1, \ldots$ such that:

- $X_t$ depends only on $X_{t-1}$
- In other words, $X_t$ is conditionally independent of $X_0, \ldots, X_{t-2}$ given $X_{t-1}$.

The evolution of the chain is governed by:

- A transition matrix $M$ (discrete case)
- A transition kernel $K$ (continuous case)

A discrete-time Markov chain is a sequence of random variables $X_0, X_1, \ldots$ such that:

- ▶ $X_t$ depends only on $X_{t-1}$
- ▶ In other words, $X_t$ is conditionally independent of $X_0, \ldots, X_{t-2}$ given $X_{t-1}$.

The evolution of the chain is governed by:

- ▶ A transition matrix $M$ (discrete case)
- ▶ A transition kernel $K$ (continuous case)
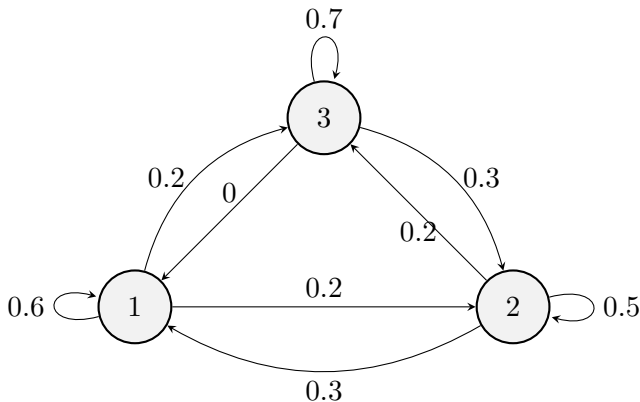
Let us denote our state-space with X.

Consider the transition matrix:

$$M = \begin{bmatrix} 0.6 & 0.2 & 0.2 \\ 0.3 & 0.5 & 0.2 \\ 0 & 0.3 & 0.7 \end{bmatrix}, \qquad \text{where X} = \{1, 2, 3\}.$$

What is a Markov chain?

Example 1: Simulate a discrete Markov chain – What does the matrix $M$ mean?

| $M$ | $X_t = 1$ | $X_t = 2$ | $X_t = 3$ |
|-----------|------|------|------|
| $X_{t-1} = 1$ | 0.6 | 0.2 | 0.2 |
| $X_{t-1} = 2$ | 0.3 | 0.5 | 0.2 |
| $X_{t-1} = 3$ | 0 | 0.3 | 0.7 |

Example: Given $X_0 = 1$, how to simulate this chain?

Sample:

$$X_t | X_t = x_{t-1} \sim \text{Discrete}(M_{x_{t-1}, \cdot}).$$

# What is a Markov chain?

Example 1: Simulate a discrete Markov chain – What does the matrix $M$ mean?

| $M$ | $X_t = 1$ | $X_t = 2$ | $X_t = 3$ |
|---|---|---|---|
| $X_{t-1} = 1$ | 0.6 | 0.2 | 0.2 |
| $X_{t-1} = 2$ | 0.3 | 0.5 | 0.2 |
| $X_{t-1} = 3$ | 0 | 0.3 | 0.7 |

Example: Given $X_0 = 1$, how to simulate this chain?

Sample:

$$X_t | X_t = x_{t-1} \sim \text{Discrete}(M_{x_{t-1},\cdot}).$$

Simulation!

Let $p_0(i) = \mathbb{P}(X_0 = i)$ for $i \in \mathrm{X}$.

Let $p_0(i) = \mathbb{P}(X_0 = i)$ for $i \in X$. Then, the density of the chain at time $n$ is given by:

$$
\begin{aligned}
p_n(i) &= \mathbb{P}(X_n = i) \\
&= \sum_k \mathbb{P}(X_n = i, X_{n-1} = k) \\
&= \sum_k \mathbb{P}(X_n = i | X_{n-1} = k) \mathbb{P}(X_{n-1} = k) \\
&= \sum_k M_{ki} p_{n-1}(k).
\end{aligned}
$$

Let $p_0(i) = \mathbb{P}(X_0 = i)$ for $i \in X$. Then, the density of the chain at time $n$ is given by:

$$
\begin{aligned}
p_n(i) &= \mathbb{P}(X_n = i) \\
&= \sum_k \mathbb{P}(X_n = i, X_{n-1} = k) \\
&= \sum_k \mathbb{P}(X_n = i | X_{n-1} = k)\mathbb{P}(X_{n-1} = k) \\
&= \sum_k M_{ki} p_{n-1}(k).
\end{aligned}
$$

This implies that

$$
p_n = p_{n-1} M.
$$

Let $p_0(i) = \mathbb{P}(X_0 = i)$ for $i \in X$. Then, the density of the chain at time $n$ is given by:

$$\begin{aligned}
p_n(i) &= \mathbb{P}(X_n = i) \\
&= \sum_k \mathbb{P}(X_n = i, X_{n-1} = k) \\
&= \sum_k \mathbb{P}(X_n = i | X_{n-1} = k) \mathbb{P}(X_{n-1} = k) \\
&= \sum_k M_{ki} p_{n-1}(k).
\end{aligned}$$

This implies that

$$p_n = p_{n-1} M.$$

Therefore,

$$p_n = p_0 M^n.$$

We need Markov chains

We need Markov chains

▶ With invariant distributions

▶ Their convergence is ensured

▶ Their invariant distribution is unique

We need Markov chains

▶ With invariant distributions

▶ Their convergence is ensured

▶ Their invariant distribution is unique

We will now review the properties which ensure these in discrete space case.

For two states, $x, x' \in X$, we write $x \rightsquigarrow x'$ if there is a path from $x$ to $x'$:

$$\exists n > 0, \text{ s.t. }, \mathbb{P}(X_n = x' | X_0 = x) > 0.$$

For two states, $x, x' \in X$, we write $x \rightsquigarrow x'$ if there is a path from $x$ to $x'$:

$$\exists n > 0, \text{ s.t. }, \mathbb{P}(X_n = x' | X_0 = x) > 0.$$

If $x \rightsquigarrow x'$ and $x' \rightsquigarrow x$, then we say that $x$ and $x'$ *communicate*.

For two states, $x, x' \in X$, we write $x \rightsquigarrow x'$ if there is a path from $x$ to $x'$:

$$\exists n > 0, \text{ s.t. }, \mathbb{P}(X_n = x' | X_0 = x) > 0.$$

If $x \rightsquigarrow x'$ and $x' \rightsquigarrow x$, then we say that $x$ and $x'$ *communicate*.

A communication class $C \subset X$ is a set of states such that $x \in C$ and $x' \in C$ if and only if $x \rightsquigarrow x'$ and $x' \rightsquigarrow x$.

For two states, $x, x' \in X$, we write $x \rightsquigarrow x'$ if there is a path from $x$ to $x'$:

$$\exists n > 0, \text{ s.t. }, \mathbb{P}(X_n = x' | X_0 = x) > 0.$$

If $x \rightsquigarrow x'$ and $x' \rightsquigarrow x$, then we say that $x$ and $x'$ *communicate*.

A communication class $C \subset X$ is a set of states such that $x \in C$ and $x' \in C$ if and only if $x \rightsquigarrow x'$ and $x' \rightsquigarrow x$.

A chain is irreducible if X is a single communication class.

A Markov chain is recurrent if every state is to be visited infinitely often.

A Markov chain is recurrent if every state is to be visited infinitely often.

Define the return time:

$$\tau_i = \inf\{n \geq 0 : X_n = i\}.$$

We say that the state is recurrent if

$$\mathbb{P}(\tau_i < \infty | X_1 = i) = 1.$$

If a state is not recurrent, it is transient.

We say that a state $i$ is positively recurrent if

$$\mathbb{E}[\tau_i | X_1 = i] < \infty.$$

We say that a state $i$ is positively recurrent if

$$\mathbb{E}[\tau_i|X_1 = i] < \infty.$$

If a recurrent state is not positive recurrent, it is null recurrent.

A probability mass function $\pi$ is called $M$-invariant if

$$\pi(i) = \sum_j M_{ij}\pi(j).$$

A probability mass function $\pi$ is called $M$-invariant if

$$\pi(i) = \sum_j M_{ij}\pi(j).$$

Equivalently

$$\pi = \pi M.$$

### Theorem 1

*If M is irreducible, then M has a unique invariant distribution if and only if it is positive recurrent.*

### Theorem 1

*If M is irreducible, then M has a unique invariant distribution if and only if it is positive recurrent.*

This is existence, we do not talk about convergence yet.

We will define reversibility through detailed balance condition.

We will define reversibility through detailed balance condition.

A Markov transition matrix $M$ is reversible w.r.t. $\pi$ if and only if for all $i, j \in X$,

$$\pi(i)M_{ij} = \pi(j)M_{ji}.$$

We will define reversibility through detailed balance condition.

A Markov transition matrix $M$ is reversible w.r.t. $\pi$ if and only if for all $i, j \in \mathrm{X}$,

$$\pi(i)M_{ij} = \pi(j)M_{ji}.$$

This is called the detailed balance condition (we will discuss the continuous version)

Constructing a chain with stationary distribution $\pi$ is ensured if detailed balance is satisfied since it implies $\pi = \pi M$.

We have seen how to construct chains with invariant distributions $\pi$.

We have seen how to construct chains with invariant distributions $\pi$.

However, the convergence of the chain $p_n \to \pi$ requires one more ingredient: ergodicity.

We have seen how to construct chains with invariant distributions $\pi$.

However, the convergence of the chain $p_n \to \pi$ requires one more ingredient: ergodicity.

For this, we need a final ingredient: aperiodicity.

A state $i$ is called aperiodic if

$$\{n > 0 : \mathbb{P}(X_{n+1} = i | X_1 = i) > 0\}$$

has no common divisor other than 1.

### Definition 2

An irreducible Markov chain is called ergodic if it is positive recurrent and aperiodic.

### Definition 2

An irreducible Markov chain is called ergodic if it is positive recurrent and aperiodic.

Ergodicity brings us the missing ingredient for the convergence: We can now ensure $p_n$ to converge to $\pi$.

If $(X_n)_{n \in \mathbb{N}}$ is an ergodic Markov chain with any initial $p_0$ and a Markov transition matrix $M$ with $\pi$ as its invariant distribution, then

$$\lim_{n \to \infty} p_n(i) = \pi(i).$$

If $(X_n)_{n \in \mathbb{N}}$ is an ergodic Markov chain with any initial $p_0$ and a Markov transition matrix $M$ with $\pi$ as its invariant distribution, then

$$\lim_{n \to \infty} p_n(i) = \pi(i).$$

Moreover, for $i, j \in \mathrm{X}$

$$\lim_{n \to \infty} \mathbb{P}(X_n = i | X_1 = j) = \pi(i).$$

What about continuous state-space Markov chains, i.e., where X is uncountable, e.g., $X = \mathbb{R}$?

What about continuous state-space Markov chains, i.e., where X is uncountable, e.g., $X = \mathbb{R}$?

We will be mainly interested in the continuous case, however, the analogous concepts are defined in a much more complicated way.

What about continuous state-space Markov chains, i.e., where X is uncountable, e.g., $X = \mathbb{R}$?

We will be mainly interested in the continuous case, however, the analogous concepts are defined in a much more complicated way.

We will not go into the details here, we will just now introduce the continuous state-space notation.

We assume now our state-space is uncountable, e.g., $X = \mathbb{R}$.

We assume now our state-space is uncountable, e.g., $X = \mathbb{R}$.

We denote the initial *density* of the chain by $p_0(x)$.

We assume now our state-space is uncountable, e.g., $X = \mathbb{R}$.

We denote the initial *density* of the chain by $p_0(x)$.

The transition kernel is denoted $K(x_n | x_{n-1})$.

We assume now our state-space is uncountable, e.g., $X = \mathbb{R}$.

We denote the initial *density* of the chain by $p_0(x)$.

The transition kernel is denoted $K(x_n|x_{n-1})$.

The density of the chain at time $n$ is denoted by $p_n(x_n)$.

A discrete-time Markov chain is a process $(X_n)_{n\in\mathbb{N}}$, when X is uncountable, satisfies:

A discrete-time Markov chain is a process $(X_n)_{n \in \mathbb{N}}$, when X is uncountable, satisfies:

$$p(x_n|x_{1:n-1}) = p(x_n|x_{n-1}) = K(x_n|x_{n-1}).$$

A discrete-time Markov chain is a process $(X_n)_{n\in\mathbb{N}}$, when X is uncountable, satisfies:

$$p(x_n|x_{1:n-1}) = p(x_n|x_{n-1}) = K(x_n|x_{n-1}).$$

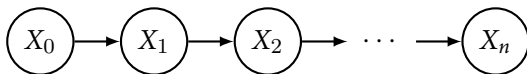We will again consider the time-homogeneous case, i.e. the transition kernel is time-independent.

We will again consider the time-homogeneous case, i.e. the transition kernel is time-independent. A Markov chain therefore can be defined entirely by its:

▶ Initial state (or initial distribution)

▶ Transition kernel

The transition kernel is a density function $K(x_n|x_{n-1})$ for fixed $x_{n-1}$, i.e.,

$$\int_X K(x_n|x_{n-1})\,\mathrm{d}x_n = 1.$$

Otherwise, it is a function of $(x_n, x_{n-1})$.

Consider the following Markov chain: $X_0 = 0$ and

$$K(x_n|x_{n-1}) = \mathcal{N}(x_n; ax_{n-1}, 1),$$

where $0 < a < 1$.

Consider the following Markov chain: $X_0 = 0$ and

$$K(x_n|x_{n-1}) = \mathcal{N}(x_n; ax_{n-1}, 1),$$

where $0 < a < 1$.
We can simulate this chain by:

$$X_1 \sim \mathcal{N}(0, 1)$$
$$X_2 \sim \mathcal{N}(aX_1, 1)$$
$$X_3 \sim \mathcal{N}(aX_2, 1)$$
$$\vdots$$
$$X_n \sim \mathcal{N}(aX_{n-1}, 1).$$

Simulation.

The Chapman-Kolmogorov equation for the continuous case

$$p(x_n|x_{n-k}) = \int_X K(x_n|x_{n-1})p(x_{n-1}|x_{n-k}) \, \mathrm{d}x_{n-1},$$

for $k > 1$.

Let $p_0(x)$ be the initial density such that $X_0 \sim p_0(x)$.

Then, the density of the chain at time $n$ is given by

$$p_n(x_n) = \int_X K(x_n|x_{n-1})p_{n-1}(x_{n-1})\,\mathrm{d}x_{n-1}.$$

It is useful for us to define the $m$-step transition kernel:

$$p(x_{m+n}|x_n) = K^m(x_{m+n}|x_n),$$
$$= \int_X K(x_{m+n}|x_{m+n-1}) \cdots K(x_{n+1}|x_n) \, dx_{m+n-1} \cdots dx_{n+1}.$$

We have the similar conditions of aperiodicity and irreducibility as in the discrete case, but,

▶ These are defined over *sets* rather than states.
▶ irreducibility is replaced by $\phi$-irreducibility.
▶ aperiodicity is defined for sets

We have the similar conditions of aperiodicity and irreducibility as in the discrete case, but,

▶ These are defined over *sets* rather than states.

▶ irreducibility is replaced by $\phi$-irreducibility.

▶ aperiodicity is defined for sets

We will not go into the details of these conditions for continuous space case.

A probability distribution $\pi$ is called $K$-invariant if

$$\pi(x) = \int_X \pi(x')K(x|x')\,\mathrm{d}x'.$$

Similar to the discrete case.

The detailed balance condition for the continuous case takes a similar form:

$$\pi(x)K(x'|x) = \pi(x')K(x|x').$$

The detailed balance condition for the continuous case takes a similar form:

$$\pi(x)K(x'|x) = \pi(x')K(x|x').$$

Note that this is a sufficient condition for stationarity of $\pi$:

$$\int \pi(x)K(x'|x)\mathrm{d}y = \int \pi(x')K(x|x')\mathrm{d}x',$$

$$\implies \pi(x) = \int K(x|x')\pi(x')\mathrm{d}x',$$

which implies $\pi$ is $K$-invariant.

A useful formulation of reversibility is the following: A Markov kernel $K$ is $\pi$-reversible if

$$\int \int f(x, x')\pi(x)K(x|x')\mathrm{d}x\mathrm{d}x' = \int \int f(x, x')\pi(x')K(x'|x)\mathrm{d}x\mathrm{d}x',$$

for every measurable $f$, which follows from the detailed balance condition.

Consider the following Markov chain: $X_0 = 0$ and

$$K(x_n|x_{n-1}) = \mathcal{N}(x_n; ax_{n-1}, 1),$$

where $0 < a < 1$.

Consider the following Markov chain: $X_0 = 0$ and

$$K(x_n|x_{n-1}) = \mathcal{N}(x_n; ax_{n-1}, 1),$$

where $0 < a < 1$. Note that we can also write this as

$$X_n = aX_{n-1} + \epsilon_n,$$

where $\epsilon_n \sim \mathcal{N}(0, 1)$.

Prove that for

$$\pi(x) = \mathcal{N}\left(x; 0, \frac{1}{1-a^2}\right),$$

the detailed balance condition is satisfied for the kernel

$$K(x_n|x_{n-1}) = \mathcal{N}(x_n; ax_{n-1}, 1),$$

where $0 < a < 1$.

Prove that $K^m(x_{m+n}|x_n)$ is given by

$$K^m(x_{m+n}|x_n) = \mathcal{N}\left(x_{m+n}; a^m x_n, \frac{1 - a^{2m}}{1 - a^2}\right).$$

Then prove that

$$\pi(x) = \lim_{m \to \infty} K^m(x|x'),$$

independent of $x'$.

### Theorem 3

*If K is an irreducible, $\pi$-invariant kernel, then for any integrable function $\varphi$*

$$\lim_{T \to \infty} \frac{1}{T} \sum_{i=1}^{T} \varphi(X_i) = \int \varphi(x)\pi(x)\mathrm{d}x = (\varphi, \pi),$$

*almost surely, for almost all initial points $x_0$.*

### Theorem 3

*If K is an irreducible, $\pi$-invariant kernel, then for any integrable function*
*$\varphi$*

$$\lim_{T \to \infty} \frac{1}{T} \sum_{i=1}^{T} \varphi(X_i) = \int \varphi(x)\pi(x)\mathrm{d}x = (\varphi, \pi),$$

*almost surely, for almost all initial points $x_0$.*

Therefore, we can use these samples to estimate our integrals.

### Theorem 4

*If $K$ is irreducible, aperiodic, and $\pi$-invariant, then*

$$\lim_{T \to \infty} \int_{X} |K^T(y|x) - \pi(y)| \mathrm{d}y = 0,$$

*for $\pi$-almost all starting values $x$.*

We will first look at a surprisingly simple approach: the Metropolis-Hastings algorithm.

We will first look at a surprisingly simple approach: the Metropolis-Hastings algorithm.

This approach relies on the following idea:

We will first look at a surprisingly simple approach: the Metropolis-Hastings algorithm.

This approach relies on the following idea:

▶ We can sample from a proposal $q(x|x')$ (that is a Markov kernel)

We will first look at a surprisingly simple approach: the Metropolis-Hastings algorithm.

This approach relies on the following idea:

▶ We can sample from a proposal $q(x|x')$ (that is a Markov kernel)

▶ We can use accept/reject

We will first look at a surprisingly simple approach: the Metropolis-Hastings algorithm.

This approach relies on the following idea:

► We can sample from a proposal $q(x|x')$ (that is a Markov kernel)

► We can use accept/reject

We can design the process so that the stationary distribution of the chain is the target distribution.

We will first look at a surprisingly simple approach: the Metropolis-Hastings algorithm.

This approach relies on the following idea:

▶ We can sample from a proposal $q(x|x')$ (that is a Markov kernel)

▶ We can use accept/reject

We can design the process so that the stationary distribution of the chain is the target distribution.

This is however very different from the rejection sampling approach.

Consider the following method:

▶ Sample $X' \sim q(x'|X_{n-1})$

▶ Set $X_n = X'$ with probability

$$\alpha(X'|X_{n-1}) = \min\left\{1, \frac{\pi(X')q(X_{n-1}|X')}{\pi(X_{n-1})q(X'|X_{n-1})}\right\}.$$

▶ Otherwise, set $X_n = X_{n-1}$.

Consider the following method:

▶ Sample $X' \sim q(x'|X_{n-1})$

▶ Set $X_n = X'$ with probability

$$\alpha(X'|X_{n-1}) = \min\left\{1, \frac{\pi(X')q(X_{n-1}|X')}{\pi(X_{n-1})q(X'|X_{n-1})}\right\}.$$

▶ Otherwise, set $X_n = X_{n-1}$.

Note the last step: we discard the sample $X'$ if rejected BUT set $X_n = X_{n-1}$.

The ratio

$$r(x, x') = \frac{\pi(x')q(x|x')}{\pi(x)q(x'|x)},$$

is called acceptance ratio.

We have discussed explicit kernels in the discrete and continuous cases.

We have discussed explicit kernels in the discrete and continuous cases.

But the MH algorithm automatically gives us a kernel.

We have discussed explicit kernels in the discrete and continuous cases.

But the MH algorithm automatically gives us a kernel.

How to prove that the stationary distribution is the target distribution?

Let us figure out the kernel.

Let us figure out the kernel.

Let us say, we have the sample from the proposal $x'$. Fixing this sample, the acceptance step samples from the mixture (*intuitively*):

$$\alpha(x'|x)\delta_{x'}(y) + (1 - \alpha(x'|x))\delta_x(y).$$

To get the full kernel, we need to integrate over $x'$:

$$K(y|x) = \int q(x'|x)\left(\alpha(x'|x)\delta_{x'}(y) + (1 - \alpha(x'|x))\delta_x(y)\right) dx',$$
$$= \alpha(y|x)q(y|x) + (1 - a(x))\delta_x(y)$$

where

$$a(x) = \int \alpha(x'|x)q(x'|x)dx'.$$

More intuition in terms of $x_n$ and $x_{n-1}$:

▶ What is the probability of being at $x_{n-1}$ and getting accepted?

$$a(x_{n-1}) = \int_X \alpha(x|x_{n-1})q(x|x_{n-1})\mathrm{d}x.$$

▶ Therefore, the probability of being at $x_{n-1}$ and getting rejected is $1 - a(x_{n-1})$.

We can see that the kernel is

$$K(x_n|x_{n-1}) = \alpha(x_n|x_{n-1})q(x_n|x_{n-1}) + (1 - a(x_{n-1}))\delta_{x_{n-1}}(x_n).$$

We can now prove that the kernel satisfies the detailed balance condition:

$$K(x'|x)\pi(x) = K(x|x')\pi(x').$$

$$\pi(x)K(x'|x) = \pi(x)q(x'|x)\alpha(x',x) + \pi(x)(1-a(x))\delta_x(x')$$
$$= \pi(x)q(x'|x)\min\left\{1, \frac{\pi(x')q(x|x')}{\pi(x)q(x'|x)}\right\} + \pi(x)(1-a(x))\delta_x(x')$$
$$= \min\left\{\pi(x)q(x'|x), \pi(x')q(x|x')\right\} + \pi(x)(1-a(x))\delta_x(x')$$
$$= \min\left\{\frac{\pi(x)q(x'|x)}{\pi(x')q(x|x')}, 1\right\}\pi(x')q(x|x') + \pi(x')(1-a(x'))\delta_{x'}(x)$$
$$= K(x|x')\pi(x').$$

Assume we are given an unnormalised density to sample $\gamma$ where

$$\pi(x) = \frac{\gamma(x)}{Z},$$

where $Z$ is the normalisation constant.

▶ Sample $X' \sim q(x'|X_{n-1})$

▶ Set $X_n = X'$ with probability

$$\alpha(X'|X_{n-1}) = \min\left\{1, \frac{\gamma(X')q(X_{n-1}|X')}{\gamma(X_{n-1})q(X'|X_{n-1})}\right\}.$$

▶ Otherwise, set $X_n = X_{n-1}$.

as the normalising constants of $\pi$ would cancel out.

- ▶ Independent proposals
- ▶ Symmetric (random walk) proposals
- ▶ Gradient-based proposals
- ▶ Adaptive proposals

Choose the proposal $q(x)$ independently of the current state $X_{n-1}$. Leads to

- $X' \sim q(x')$
- Accept with probability

$$\alpha(X'|X_{n-1}) = \min\left\{1, \frac{\pi(X')q(X_{n-1})}{\pi(X_{n-1})q(X')}\right\}.$$

- Otherwise, set $X_n = X_{n-1}$.

Let us say

$$\pi(x) = \mathcal{N}(x; \mu, \sigma^2)$$

For the example, assume we want to use MH to sample from it. Choose a proposal

$$q(x) = \mathcal{N}(x; \mu_q, \sigma_q^2).$$

How to compute the acceptance ratio?

$$r(x, x') = \frac{\pi(x')q(x)}{\pi(x)q(x')}$$

$$= \frac{\mathcal{N}(x'; \mu, \sigma^2)\mathcal{N}(x; \mu_q, \sigma_q^2)}{\mathcal{N}(x; \mu, \sigma^2)\mathcal{N}(x'; \mu_q, \sigma_q^2)}$$

$$= \frac{\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x'-\mu)^2}{2\sigma^2}\right) \frac{1}{\sqrt{2\pi\sigma_q^2}} \exp\left(-\frac{(x-\mu_q)^2}{2\sigma_q^2}\right)}{\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \frac{1}{\sqrt{2\pi\sigma_q^2}} \exp\left(-\frac{(x'-\mu_q)^2}{2\sigma_q^2}\right)}$$

$$= \frac{\exp\left(-\frac{(x'-\mu)^2}{2\sigma^2}\right) \exp\left(-\frac{(x-\mu_q)^2}{2\sigma_q^2}\right)}{\exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \exp\left(-\frac{(x'-\mu_q)^2}{2\sigma_q^2}\right)}$$

$$= e^{\left(-\frac{1}{2\sigma^2}\left[(x'-\mu)^2 - (x-\mu)^2\right]\right)} e^{\left(-\frac{1}{2\sigma_q^2}\left[(x-\mu_q)^2 - (x'-\mu_q)^2\right]\right)}$$

We can choose:

$$q(x'|x) = \mathcal{N}(x'; x, \sigma_q^2)$$

The proposal looks at where we are and take a random step (random walk).

We can choose:

$$q(x'|x) = \mathcal{N}(x'; x, \sigma_q^2)$$

The proposal looks at where we are and take a random step (random walk).

Note that $q(x'|x)$ is symmetric, i.e. $q(x|x') = q(x'|x)$.

Acceptance ratio:

$$
\begin{aligned}
\mathrm{r}(x, x') &= \frac{\pi(x')q(x|x')}{\pi(x)q(x'|x)} \\
&= \frac{\pi(x')}{\pi(x)}, \\
&= \frac{\mathcal{N}(x'; \mu, \sigma^2)}{\mathcal{N}(x; \mu, \sigma^2)} \\
&= e^{\left(-\frac{1}{2\sigma^2}\left[(x'-\mu)^2 - (x-\mu)^2\right]\right)}.
\end{aligned}
$$

Set a burnin period:

► Run the sampler for fixed number of iterations and discard the first *n* samples.

► This accounts for the convergence to the stationary measure.

We can *inform* the proposal by using the gradient of the target distribution.

$$q(x'|x) = \mathcal{N}(x'; x + \gamma \nabla \log \pi(x), 2\gamma I),$$

This tends to behave really well.

We can *inform* the proposal by using the gradient of the target distribution.

$$q(x'|x) = \mathcal{N}(x'; x + \gamma \nabla \log \pi(x), 2\gamma I),$$

This tends to behave really well.

This approach is called *Metropolis adjusted Langevin algorithm* (MALA). (more on these later)

► One has to be careful that $p/q < \infty$ (while no theoretical reason, the performance tends to be quite bad).

- ▶ One has to be careful that $p/q < \infty$ (while no theoretical reason, the performance tends to be quite bad).
- ▶ The proposal should attain a balance of acceptance rate and efficiency.

- One has to be careful that $p/q < \infty$ (while no theoretical reason, the performance tends to be quite bad).
- The proposal should attain a balance of acceptance rate and efficiency.
- Too high acceptance rate is **not** necessarily good: You might be taking too small steps and getting stuck in some regions

Let us look at now the Bayesian inference problem.

We can solve it in full generality (in theory) using MH.

Recall the general formulation

$$p(x|y_{1:n}) = \frac{p(y_{1:n}|x)p(x)}{p(y_{1:n})} = \frac{\prod_{i=1}^{n} p(y_i|x)p(x)}{p(y_{1:n})},$$

when $y_1, \ldots, y_n$ are conditionally independent given $x$.

We write

$$p(x|y_{1:n}) \propto \prod_{i=1}^{n} p(y_i|x)p(x),$$

and set

$$\gamma(x) = \prod_{i=1}^{n} p(y_i|x)p(x),$$

as our unnormalised posterior.

The generic MH for Bayesian inference, given $x_{n-1}$

- Sample $X' \sim q(x'|x_{n-1})$.
- Accept $x_n = x'$ with probability

$$\alpha(x_{n-1}, x') = \min \left\{ 1, \frac{\gamma(x')q(x_{n-1}|x')}{\gamma(x_{n-1})q(x'|x_{n-1})} \right\}.$$

- Otherwise, $X_n = x_{n-1}$.

Recall our example about localising a source using observations from a sensor network.

We can now formalise this problem. Assume that the source is located at $x \in \mathbb{R}^2$ and the sensor network is located at $s_1, \ldots, s_3 \in \mathbb{R}^2$ (3 sensors).

Assume that these three sensors "observe" the source according to:

$$p(y_i|x, s_i) = \mathcal{N}(y_i; \|x - s_i\|, R),$$

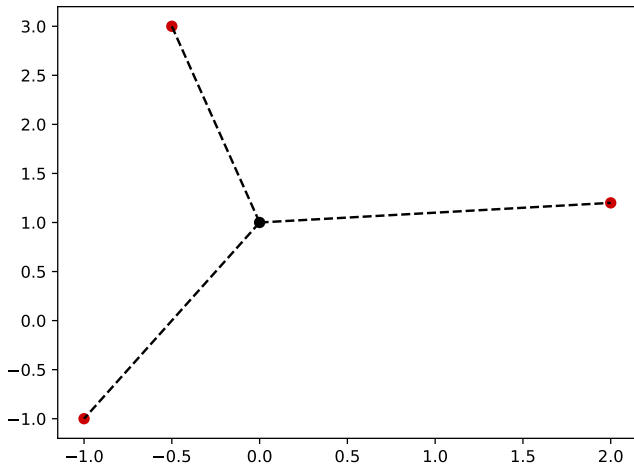where $y_i$ is the observation from sensor $i$.

Figure: Source localisation

Assume that you are asked to estimate the location of the source given
the observations $y_1, y_2, y_3$. What is the model?

Assume that you are asked to estimate the location of the source given the observations $y_1, y_2, y_3$. What is the model?

We first need a prior on the source location:

$$p(x) = \mathcal{N}(x; \mu, \Sigma),$$

where $\mu$ is the prior mean and $\Sigma$ is the prior covariance. We already have the likelihoods for each $y_i$.

The posterior is given by

$$p(x|y_1, y_2, y_3, s_1, s_2, s_3) \propto p(x) \prod_{i=1}^{3} p(y_i|x, s_i).$$

We choose a random walk proposal:

$$q(x'|x) = \mathcal{N}(x'; x, \sigma^2 I).$$

This is symmetric so the acceptance ratio is:

$$r(x, x') = \frac{p(x')p(y_1|x', s_1)p(y_2|x', s_2)p(y_3|x', s_3)}{p(x)p(y_1|x, s_1)p(y_2|x, s_2)p(y_3|x, s_3)}.$$

Consider the 2D density

$$p(x, y) \propto \exp\left(-\frac{x^2}{10} - \frac{y^4}{10} - 2(y - x^2)^2\right).$$
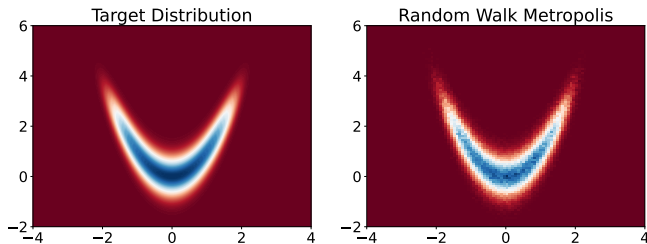
Assume we would like to sample from it.

Figure: The banana density (unnormalised)

We have

$$\gamma(x, y) = \exp\left(-\frac{x^2}{10} - \frac{y^4}{10} - 2(y - x^2)^2\right).$$

and let us choose two alternative proposals

▶ The random walk proposal:

$$q(x', y'|x, y) = \mathcal{N}(x'; x, \sigma_q^2)\mathcal{N}(y'; y, \sigma_q^2).$$

▶ and the gradient-based proposal (MALA):

$$q(x', y'|x, y) = \mathcal{N}(z; z + \gamma \nabla \log \gamma(z), \sqrt{2\gamma}\mathbf{I}).$$

where $z = (x, y)$ and $\gamma$ is a step size.

We have seen Metropolis-Hastings sampler.

We have seen Metropolis-Hastings sampler.

▶ Unfortunately, it may not be very efficient.

We have seen Metropolis-Hastings sampler.

▶ Unfortunately, it may not be very efficient.

▶ Acceptance ratios are very tricky to compute in a variety of settings:

We have seen Metropolis-Hastings sampler.

- ▶ Unfortunately, it may not be very efficient.
- ▶ Acceptance ratios are very tricky to compute in a variety of settings:
    - ▶ High-dimensional problems

We have seen Metropolis-Hastings sampler.

- ▶ Unfortunately, it may not be very efficient.
- ▶ Acceptance ratios are very tricky to compute in a variety of settings:
  - ▶ High-dimensional problems
  - ▶ Complex models

We have seen Metropolis-Hastings sampler.

▶ Unfortunately, it may not be very efficient.
▶ Acceptance ratios are very tricky to compute in a variety of settings:
  ▶ High-dimensional problems
  ▶ Complex models
  ▶ Large datasets

We have seen Metropolis-Hastings sampler.

- ▶ Unfortunately, it may not be very efficient.
- ▶ Acceptance ratios are very tricky to compute in a variety of settings:
  - ▶ High-dimensional problems
  - ▶ Complex models
  - ▶ Large datasets
- ▶ We will now look at a different approach: Langevin MCMC.

Consider the Langevin SDE for a generic drift $\nabla V$:

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

where $(B_t)_{t \geq 0}$ is a Brownian motion.

Consider the Langevin SDE for a generic drift $\nabla V$:

$$\mathrm{d}X_t = -\nabla V(X_t)\mathrm{d}t + \sqrt{2}\mathrm{d}B_t,$$

where $(B_t)_{t\geq 0}$ is a Brownian motion. This SDE has a stationary measure

$$\pi \propto e^{-V(x)}.$$

Therefore, for a classical *sampling* problem for, say $\pi(x)$, we could set
$V(x) = -\log \pi(x)$ (negative density).

Consider the Langevin SDE for a generic drift $\nabla V$:

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

where $(B_t)_{t\geq 0}$ is a Brownian motion. This SDE has a stationary measure

$$\pi \propto e^{-V(x)}.$$

Therefore, for a classical *sampling* problem for, say $\pi(x)$, we could set $V(x) = -\log \pi(x)$ (negative density).

This diffusion converges to its stationary measure exponentially fast if $V$ is $\mu$-strongly-convex.

Consider the Langevin SDE for a generic drift $\nabla V$:

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

where $(B_t)_{t \geq 0}$ is a Brownian motion.

Consider the Langevin SDE for a generic drift $\nabla V$:

$$\mathrm{d}X_t = -\nabla V(X_t)\mathrm{d}t + \sqrt{\frac{2}{\beta}}\mathrm{d}B_t,$$

where $(B_t)_{t\geq 0}$ is a Brownian motion. This SDE has a stationary measure

$$\pi \propto e^{-\beta V(x)}.$$

Consider the Langevin SDE for a generic drift $\nabla V$:

$$\mathrm{d}X_t = -\nabla V(X_t)\mathrm{d}t + \sqrt{\frac{2}{\beta}}\mathrm{d}B_t,$$

where $(B_t)_{t \geq 0}$ is a Brownian motion. This SDE has a stationary measure

$$\pi \propto e^{-\beta V(x)}.$$

This stationary measure concentrates on the minima of $V$ as $\beta \to \infty$ (Hwang, 1980).

Consider the Langevin SDE for a generic drift $\nabla V$:

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

where $(B_t)_{t \geq 0}$ is a Brownian motion. This SDE has a stationary measure

$$\pi \propto e^{-\beta V(x)}.$$

This stationary measure concentrates on the minima of $V$ as $\beta \to \infty$ (Hwang, 1980).

Langevin diffusion is a global optimiser.

The Euler discretisation is the *unadjusted Langevin algorithm* (ULA):

$$X_{t+1}^{\gamma} = X_t^{\gamma} - \gamma \nabla V(X_t^{\gamma}) + \sqrt{2\gamma} W_{t+1}$$

where $(W_t)_{t \geq 0}$ are i.i.d standard Normal random variables.

# Langevin-based approaches

Crash course on Langevin SDE - III: Numerical discretisation

The Euler discretisation is the *unadjusted Langevin algorithm* (ULA):

$$X_{t+1}^{\gamma} = X_t^{\gamma} - \gamma \nabla V(X_t^{\gamma}) + \sqrt{2\gamma} W_{t+1}$$

where $(W_t)_{t \geq 0}$ are i.i.d standard Normal random variables.

This chain has a *different* stationary measure $\pi^{\gamma}$ but a number of guarantees can be derived for its convergence.

## Theorem 1 (Durmus and Moulines, 2019)

*Let $\mathcal{L}(X_t)$ be the law of the iterates of ULA, then*

$$W_2^2(\mathcal{L}(X_t^{\gamma}), \pi) \lesssim \left(1 - \frac{\gamma \kappa}{2}\right)^{t+1} (d/m + \|x - x^\star\|^2) + \gamma,$$

*under suitable regularity conditions for $V$, restriction on $\gamma$ where $\kappa := \kappa(m, L)$.*

An important note here is that, we can sample from the posterior $p(x|y)$ using ULA as

$$p(x|y) \propto p(x, y),$$

and

$$X_{n+1}^\gamma = X_n^\gamma + \gamma \nabla \log p(X_n^\gamma, y) + \sqrt{2}\gamma W_{n+1}.$$

An important note here is that, we can sample from the posterior $p(x|y)$ using ULA as

$$p(x|y) \propto p(x, y),$$

and

$$X_{n+1}^{\gamma} = X_n^{\gamma} + \gamma \nabla \log p(X_n^{\gamma}, y) + \sqrt{2}\gamma W_{n+1}.$$

We can see that this algorithm would approximately sample from $p(x|y)$.

Let us say we have data $y_1, \ldots, y_M$ for $M$ large. We can write the posterior as

$$p(x|y_{1:M}) \propto p(x) \prod_{i=1}^{M} p(y_i|x).$$

therefore, our potential becomes

$$V(x) = -\log p(x) - \sum_{i=1}^{M} \log p(y_i|x).$$

Let us say we have data $y_1, \ldots, y_M$ for $M$ large. We can write the posterior as

$$p(x|y_{1:M}) \propto p(x) \prod_{i=1}^{M} p(y_i|x).$$

therefore, our potential becomes

$$V(x) = -\log p(x) - \sum_{i=1}^{M} \log p(y_i|x).$$

Mini-quiz: What is the problem with MALA (or MH in general) in this case?

A similar problem of course would be for ULA.

A similar problem of course would be for ULA.

However, we can resolve this, as we can approximate the gradient using subsampling:

$$\nabla V(x) = \nabla \log p(x) + \sum_{i=1}^{M} \nabla \log p(y_i|x),$$

$$\approx \nabla \log p(x) + \frac{M}{m} \sum_{j=1}^{m} \nabla \log p(y_{k_j}|x) = \widehat{\nabla V(x)},$$

where $k_j \sim \text{Unif}\{1, \ldots, M\}$, for $j = 1, \ldots, m$ for $m \ll M$.

A similar problem of course would be for ULA.

However, we can resolve this, as we can approximate the gradient using subsampling:

$$\nabla V(x) = \nabla \log p(x) + \sum_{i=1}^{M} \nabla \log p(y_i|x),$$

$$\approx \nabla \log p(x) + \frac{M}{m} \sum_{j=1}^{m} \nabla \log p(y_{k_j}|x) = \widehat{\nabla V(x)},$$

where $k_j \sim \mathrm{Unif}\{1, \dots, M\}$, for $j = 1, \dots, m$ for $m \ll M$.

Stochastic gradients.

One can run ULA with stochastic gradients:

$$X_{n+1}^{\gamma} = X_n^{\gamma} - \gamma \widehat{\nabla V(X_n^{\gamma})} + \sqrt{2}\gamma W_{n+1}.$$

One can run ULA with stochastic gradients:

$$X_{n+1}^\gamma = X_n^\gamma - \gamma \widehat{\nabla V(X_n^\gamma)} + \sqrt{2}\gamma W_{n+1}.$$

The resulting method is called *stochastic gradient Langevin dynamics* (SGLD) (Welling and Teh, 2011).

One can run ULA with stochastic gradients:

$$X_{n+1}^\gamma = X_n^\gamma - \gamma \widehat{\nabla V(X_n^\gamma)} + \sqrt{2}\gamma W_{n+1}.$$

The resulting method is called *stochastic gradient Langevin dynamics* (SGLD) (Welling and Teh, 2011).

▶ Widely used for large-scale datasets.

One can run ULA with stochastic gradients:

$$X_{n+1}^{\gamma} = X_n^{\gamma} - \gamma \widehat{\nabla V(X_n^{\gamma})} + \sqrt{2}\gamma W_{n+1}.$$

The resulting method is called *stochastic gradient Langevin dynamics* (SGLD) (Welling and Teh, 2011).

▶ Widely used for large-scale datasets.

▶ It has similar guarantees to ULA in Wasserstein-2 distance for strongly convex $V$.

One can run ULA with stochastic gradients:

$$X_{n+1}^\gamma = X_n^\gamma - \gamma \widehat{\nabla V(X_n^\gamma)} + \sqrt{2}\gamma W_{n+1}.$$

The resulting method is called *stochastic gradient Langevin dynamics* (SGLD) (Welling and Teh, 2011).

▶ Widely used for large-scale datasets.

▶ It has similar guarantees to ULA in Wasserstein-2 distance for strongly convex $V$.

▶ Also used to model and analyse the behaviour of stochastic gradient descent methods (SGD) in deep learning.

One can run ULA with stochastic gradients:

$$X_{n+1}^{\gamma} = X_n^{\gamma} - \gamma \widehat{\nabla V(X_n^{\gamma})} + \sqrt{2}\gamma W_{n+1}.$$

The resulting method is called *stochastic gradient Langevin dynamics* (SGLD) (Welling and Teh, 2011).

▶ Widely used for large-scale datasets.

▶ It has similar guarantees to ULA in Wasserstein-2 distance for strongly convex $V$.

▶ Also used to model and analyse the behaviour of stochastic gradient descent methods (SGD) in deep learning.

Web based simulations if time permits.

📎 Durmus, Alain and Eric Moulines (2019). "High-dimensional Bayesian inference via the unadjusted Langevin algorithm". In: *Bernoulli* 25.4A, pp. 2854–2882.

📎 Hwang, Chii-Ruey (1980). "Laplace's method revisited: weak convergence of probability measures". In: *The Annals of Probability*, pp. 1177–1182.

📎 Li, Bo, Thomas Bengtsson, and Peter Bickel (2005). "Curse-of-dimensionalit revisited: Collapse of importance sampling in very large scale systems". In: *Rapport technique* 85, p. 205.

📎 Welling, Max and Yee W Teh (2011). "Bayesian learning via stochastic gradient Langevin dynamics". In: *Proceedings of the 28th international conference on machine learning (ICML-11)*. Citeseer, pp. 681–688.