

## Sampling Techniques - Continued.

02/02/2025

### \* Inverse Transform Sampling.

Method for generating random samples from a Probability distribution  $P(x)$  using cumulative distribution function (CDF). It relies on probability integral transform & its inverse.

$X \rightarrow$  R.V. with a continuous CDF  $F_X(x)$ ,  
then R.V.  $Y = F_X(x) \rightarrow$  uniformly distributed  
on  $[0, 1]$ .

$$Y = F_X(x) \sim U[0, 1]$$

CDF  $F_X(x)$  maps values of  $x$  to prob. in  $[0, 1]$ .

$Y \sim U[0,1]$ , R.V.  $\hat{X} = F_X^{-1}(Y)$  has the same dist as  $X$ :-

$$\hat{X} = F_X^{-1}(Y) \sim X$$

Inverse CDF  $F_X^{-1}(y)$  maps prob.  $y \in [0,1)$  back to values of  $X$ .

Applying  $F_X^{-1}$  to uniformly distributed samples  $Y$ , we obtain  $\hat{X}$  that follows the dist. of  $X$ .

### Algorithm

- Compute CDF  $F_X(x)$

eg:-  $P(x)$  - target Dist, then. CDF

$$F_X(x) = \int_{-\infty}^x P(t) dt.$$

- Compute the inverse CDF  $F_X^{-1}(y)$

-  $F_X^{-1}(y)$  maps  $y \in [0,1)$  to  $x$ .

- Generate uniform samples  $Y \sim U[0,1]$   
↓  
any random no. generator from computer.

- Transform  $Y$  to  $\hat{X}$ .

- Apply the inverse CDF to the uniform samples.

$$\hat{X} = F_X^{-1}(Y).$$

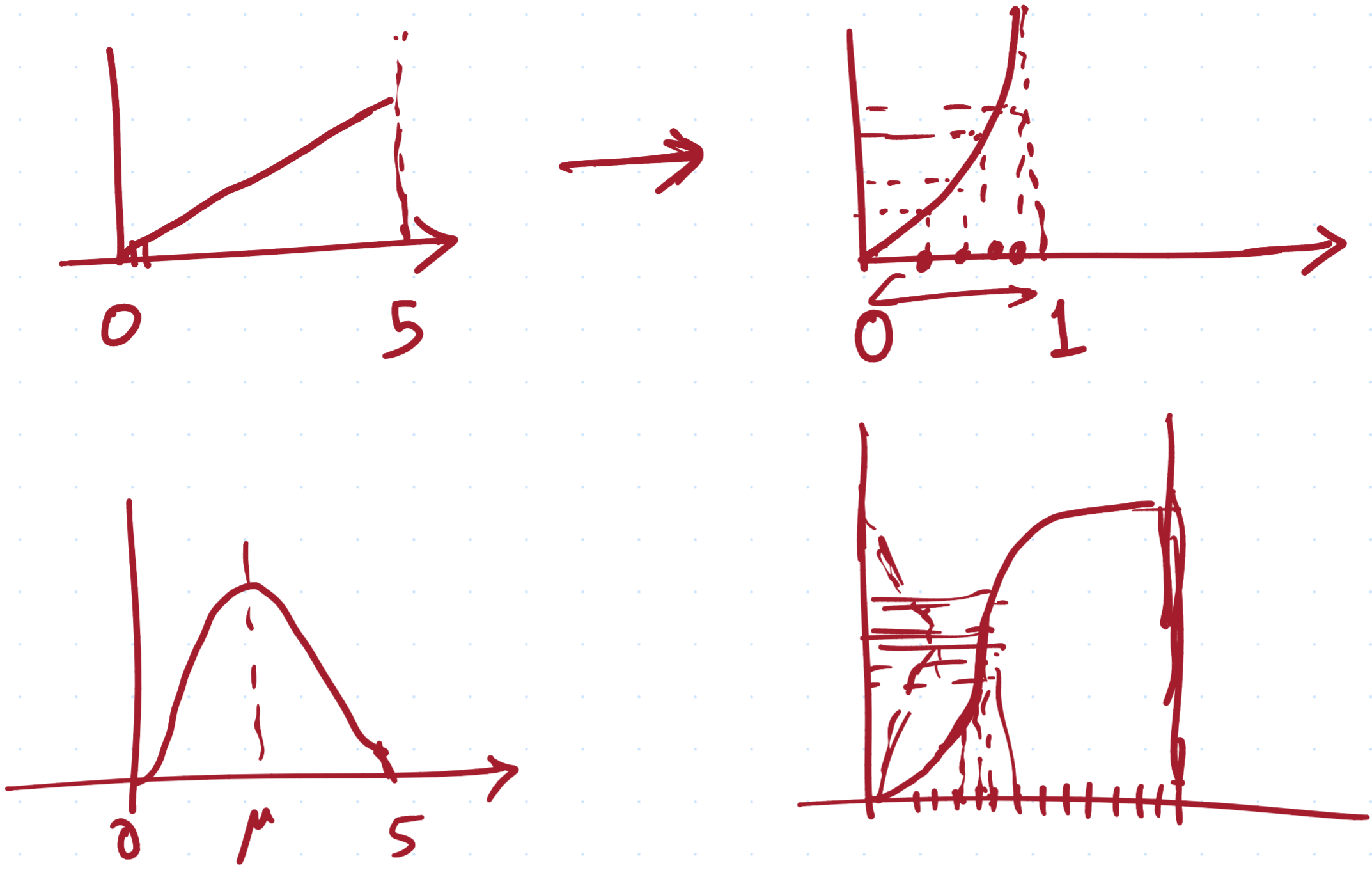
-  $\hat{X}$  will produce the target dist.  $P(X)$ .

\* Works with any dist. with a computable inverse CDF.  $F_X^{-1}(y)$

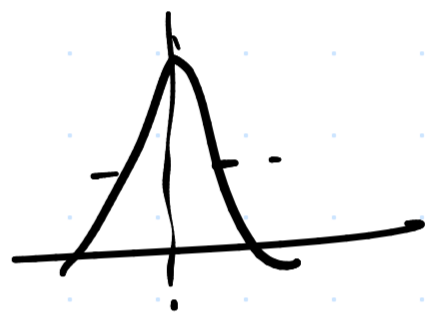
\* Computationally efficient,  $F_X^{-1}(y)$  - easy to compute

\* Dists without closed form inverse CDF,  
we use numerical methods

\*  $F_X(x) \rightarrow$  uniformly dist,  $F_X^{-1}(y)$  transforms uniform samples into samples from  $P(X)$ .



# Cauchy Distribution



PDF: Cauchy  $(x_0, \gamma) =$   
 peak of the dist  $\leftarrow x_0 \rightarrow$  location parameter  
 half width at half max  $\rightarrow \gamma > 0 \rightarrow$  scale parameter

$$\frac{1}{\pi \gamma \left[ 1 + \left( \frac{x - x_0}{\gamma} \right)^2 \right]}$$

CDF:

$$P(x \leq x) = \frac{1}{2} + \frac{1}{\pi} \tan^{-1} \left( \frac{x - x_0}{\gamma} \right)$$

Cauchy (0,1)

$$F_x^{-1}(y) = \tan \left[ \pi \left( y - \frac{1}{2} \right) \right]$$

$$F_X(x) = \int_{-\infty}^x \frac{1}{\pi\gamma \left[ 1 + \left( \frac{t-x_0}{\gamma} \right)^2 \right]} dt$$

$$u = \frac{t-x_0}{\gamma} \quad du = \frac{1}{\gamma} dt, \quad dt = \gamma du.$$

$$t \rightarrow -\infty, \quad u \rightarrow -\infty$$

$$t \rightarrow x, \quad u \rightarrow \frac{x-x_0}{\gamma}$$

$$\int_{-\infty}^{\frac{x-x_0}{\gamma}} \frac{1}{\pi\gamma [1+u^2]} \gamma du$$

$$\Rightarrow \frac{1}{\pi} \int_{-\infty}^{\frac{x-x_0}{\gamma}} \frac{1}{1+u^2} du = \frac{1}{\pi} \left[ \tan^{-1} u \right]_{-\infty}^{\frac{x-x_0}{\gamma}}$$

$$= \frac{1}{\pi} \left[ \tan^{-1} \left( \frac{x-x_0}{\gamma} \right) - \left( -\frac{\pi}{2} \right) \right]$$

as  $u \rightarrow -\infty$   
 $\tan^{-1} u \rightarrow -\frac{\pi}{2}$

$u \rightarrow \frac{x-x_0}{\gamma}$

$\tan^{-1} u \rightarrow \tan^{-1} \left( \frac{x-x_0}{\gamma} \right)$

$$= \frac{1}{\pi} \left[ \frac{\pi}{2} + \tan^{-1} \left( \frac{x-x_0}{\gamma} \right) \right]$$

$$F_X(x) = \frac{1}{2} + \frac{1}{\pi} \tan^{-1} \left( \frac{x-x_0}{\gamma} \right)$$

Cauchy  $(0, 1)$   
 $\downarrow$   $\downarrow$   
 $x$   $y$

$\rightarrow$

$$F_x(x) = \frac{1}{2} + \frac{1}{\pi} \tan^{-1}(x)$$

$$y = F_x(x)$$

$$F_x^{-1}(y) = ?$$

$$\Rightarrow y = \frac{1}{2} + \frac{1}{\pi} \tan^{-1}(x)$$

$$\Rightarrow y - \frac{1}{2} = \frac{1}{\pi} \tan^{-1}(x)$$

$$\Rightarrow \tan^{-1}(x) = \left(y - \frac{1}{2}\right) \pi$$

$$\Rightarrow x = \tan\left(\pi\left(y - \frac{1}{2}\right)\right)$$

$$x = F_x^{-1}(y) = \tan\left(\pi\left(y - \frac{1}{2}\right)\right)$$

## Exponential Dist :-

The PDF is given by :-

$$\textcircled{*} \underline{f(t)} = \underline{\text{Exp}(\lambda)} = \lambda e^{-\lambda x}, \quad \forall x \in \underline{\underline{[0, \infty)}}$$

$$\textcircled{*} \underline{F_x(x)} = \underline{P(x \leq x)} = \underline{1 - e^{-\lambda x}}$$

$$\textcircled{*} \underline{F_x^{-1}(y)} = -\log_e(1-y) \rightarrow \text{for Exp(1)}$$

$$F_x(x) = \int_{-\infty}^x \underline{f(t)} dt = \int_{-\infty}^x \lambda e^{-\lambda t} dt$$

$$= \lambda \int_0^x e^{-\lambda t} dt = \lambda \left[ \frac{e^{-\lambda t}}{-\lambda} \right]_0^x$$

$$= -1 [e^{-\lambda x} - e^0]$$

$$= 1 - e^{-\lambda x} + e$$

$$F_x(x) = 1 - e^{-\lambda x} = y$$

$$y = 1 - e^{-\lambda x}$$

$$(y-1) = -e^{-\lambda x}$$

$$(1-y) = e^{-\lambda x}$$

$$\log_e(1-y) = -\lambda x$$

$$\log_e e = ? \downarrow$$
$$\ln$$

$$\lambda = \text{Exp}(1) = ?$$

$$x = -\frac{1}{\lambda} \log_e(1-y)$$

$$x = -\log_e(1-y)$$

## Gumbel Distribution

PDF

$$\text{Gumbel}(\mu, \beta) = \frac{e^{-(z+e^{-z})}}{\beta}$$

CDF

$$P(x \leq \alpha) = e^{-e^{-\frac{\alpha-\mu}{\beta}}}$$

$$z = \frac{\alpha-\mu}{\beta}$$

$\mu \rightarrow$  location parameter  
 $\beta > 0$ , scale param.

$$\text{For Gumbel}(0,1) \Rightarrow F_x^{-1}(y) = -\log_e(-\log_e(y))$$

wikipedia

$\mu$   $\beta$



$$F_X(x) = \int_{-\infty}^x P(t) dt$$

X

$$= \int_{-\infty}^x \frac{e^{-(z + e^{-z})}}{\beta} dt$$

$$= \int_{-\infty}^{\frac{x-\mu}{\beta}} \frac{e^{-(z + e^{-z})}}{\beta} dz \quad \left| \begin{array}{l} t \rightarrow \infty, z \rightarrow \infty \\ t \rightarrow x, z \rightarrow \frac{x-\mu}{\beta} \end{array} \right.$$

$$z = \frac{t-\mu}{\beta}$$

$$dz = \frac{dt}{\beta}$$

$$dt = \beta dz$$

$$F_X(x) = \int_{-\infty}^{\frac{x-\mu}{\beta}} e^{-(z + e^{-z})} dz \quad ?$$

$$\ln u = -z,$$

$$u = e^{-z} \quad z = -\ln u$$

$$du = e^{-z} (-1) dz$$

$$\boxed{dz = -\frac{1}{u} du}$$

$$\left\{ \begin{array}{l} z = -\infty, u = \frac{1}{e^z} = \frac{1}{e^{-\infty}} \\ u \rightarrow \infty \\ = e^{\infty} \\ = \infty \end{array} \right.$$

$$z \rightarrow \frac{x-\mu}{\beta}$$

$$u \rightarrow e^{-\frac{x-\mu}{\beta}}$$

$$F_x(x) = \int_{\infty}^{e^{\frac{x-\mu}{\beta}}} e^{-(-\ln(u) + u)} \left(-\frac{1}{u}\right) du$$

$$\int_{\infty}^{e^{-\frac{x-\mu}{\beta}}} (u + e^{-u}) \left(-\frac{1}{u}\right) du$$

merged up here!!



Homework

$$F_x(x) = \int_{e^{-\frac{x-\mu}{\beta}}}^{\infty} e^{-u} du = [-e^{-u}]_{e^{-\frac{x-\mu}{\beta}}}^{\infty}$$

$$u \rightarrow \infty, -e^{-u} \rightarrow 0$$

$$u = e^{-\frac{(x-\mu)}{\beta}}, -e^{-u}$$

$$\rightarrow -e^{-e^{-\frac{(x-\mu)}{\beta}}}$$

$$F_x(x) = e^{-e^{-\frac{(x-\mu)}{\beta}}}$$



$$\Rightarrow e^{-e^{-\frac{(x-\mu)}{\beta}}} = y$$

$$\Rightarrow -e^{-\frac{(x-\mu)}{\beta}} = \ln y$$

$$\Rightarrow -\frac{x-\mu}{\beta} = \ln(-\ln(y))$$

$$\Rightarrow \frac{x-\mu}{\beta} = -\ln(-\ln(y))$$

$$\Rightarrow x = \underbrace{-\beta}_{\downarrow 1} \ln(-\ln(y)) + \underbrace{\mu}_{\downarrow 0}$$

$$\mu = 0, \beta = 1$$

$$F_X^{-1}(y) =$$

$$\boxed{x = -\ln(-\ln(y))}$$



# Langevin Monte Carlo

Stochastic Gradient Langevin Dynamics (SGLD) is a Markov chain Monte Carlo (MCMC) method used to sample from a target dist.

It combines gradient based optimization and stochastic differential equations.

Goal: Sample from a target dist  $P(x)$  with probability density  $e^{-U(x)}$ ,  $U(x) \rightarrow$  potential function.

Simulates a stochastic process, Langevin SDE that explores space in a way that converges to a target dist.

$$\left\{ \underline{dX_t} = - \underline{\nabla U(X_t)} dt + \underline{\sqrt{2}} dB_t \right\}$$

$x_t \rightarrow$  state of the system at time  $t$

$U(x_t) \rightarrow$  potential function related to the target distribution.

$\nabla U(x_t) \rightarrow$  gradient of potential function -

$dB_t \rightarrow$  increment of a standard Brownian Motion (Wiener Process)

$\sqrt{2} \rightarrow$  scaling factor for the noise.

$-\nabla U(x_t)dt \rightarrow$  drift term, pushes  $x_t$  towards regions of lower potentials, higher probability under  $P(x)$ .

$\sqrt{2}dB_t \rightarrow$  diffusion term  $\rightarrow$  add random noise, to mitigate stuckness in local minima.

$P(x) \propto e^{-\underline{U(x)}}$   $\rightarrow$  stationary dist is the target dist.

After running for a long time, the samples  $x_t$  will be distributed according to  $P(x)$ .

Discretized to generate samples:-

$$\left\{ \underline{x_{t+1}} = \underline{x_t} - \eta \underline{\nabla U(x_t)} + \underline{\sqrt{2\eta} z_t} \right\}$$

$\eta \rightarrow$  step size / learning rate

$z_t \rightarrow$  standard normal R.V.  $z_t \sim \mathcal{N}(0, I)$

$-\eta \nabla U(x_t) \rightarrow$  performs a gradient step towards regions of lower potential.

$\sqrt{2\eta} z_t \rightarrow$  adds Gaussian noise to the exploration.

Algorithm :-

Start with initial point  $x_0$ .

For each step:-

- Compute gradient  $\nabla U(x_t)$
- update

$$x_{t+1} = x_t - \eta \nabla U(x_t) + \sqrt{2\eta} z_t$$

- Store  $x_{t+1}$  as a sample.

Example :- - Sampling from a Gaussian dist.

target  $\rightarrow$  Gaussian  $N(\mu, \sigma^2)$

$$\underline{U(x)} = \frac{(x-\mu)^2}{2\sigma^2}$$

$$\nabla U(x) = \frac{x-\mu}{\sigma^2}$$

$$\alpha e^{-U(x)}$$

~~$\frac{1}{\sqrt{2\pi\sigma^2}}$~~   $e^{-\frac{(x-\mu)^2}{2\sigma^2}}$



## Langvin Update :-

$$X_{t+1} = X_t - \eta \frac{X - \mu}{\sigma^2} + \sqrt{2\eta} z_t$$

---

$$U(z) = \frac{1}{2} \left( \frac{\|z\| - 2}{0.4} \right)^2 - \log \left( e^{-0.5 \left[ \frac{z_1 - 2}{0.6} \right]} + e^{-0.5 \left[ \frac{z_1 + 2}{0.6} \right]} \right)$$

$$p(z) \propto e^{-U(z)}$$

$$\|z\| = \sqrt{z_1^2 + z_2^2}$$

→ Euclidean norm of  $z$

$$\frac{1}{2} \left( \frac{\|z\| - 2}{0.4} \right)^2 \rightarrow$$

quadratic penalty that encourages  $\|z\|$  to be close to 2

$$-\log(\dots) \rightarrow$$

log-sum-exp → creates two modes,  $z_1 = 2, z_1 = -2$

Target dist =  $p(z)$

$$p(z) \propto e^{-U(z)}$$

$$p(z) \propto \exp\left(-\frac{1}{2} \left(\frac{\|z\| - 2}{0.4}\right)^2\right) \cdot \left( e^{-0.5 \left[\frac{z_1 - 2}{0.6}\right]^2} + e^{-0.5 \left[\frac{z_1 - 2}{0.6}\right]^2} \right)$$

ring shaped gaussian  
centered at  $\|z\| = 2$   
& radius = 2 (encomp)

↓  
two modes

$$z_1 = 2, z_1 = -2$$

(peaks)

symmetric abt.

$$z_1 = 0.$$

To sample from  $p(z)$ ,  
we use Langevin Monte Carlo  
or other MCMC methods.

$$\nabla U(z) = \begin{bmatrix} \frac{\partial U}{\partial z_1} \\ \frac{\partial U}{\partial z_2} \end{bmatrix} \approx$$

$P(x) \propto e^{-U(x)} \rightarrow$  Related to energy based models; states with lower energy are more likely (higher probability) while states with higher energy are less likely (lower probability).

$P(x) \propto e^{-0} \propto \textcircled{1}$ .  $\underline{U(x) \rightarrow 0}$  lower energy

$P(x) \propto e^{-\infty} \propto \frac{1}{\infty} \propto 0$ .  $U(x) \rightarrow \infty \rightarrow$  less likely higher energy

$P(x) =$  always non-ve, the dist is normalized.

$$P(x) = \frac{1}{Z} e^{-U(x)} \quad Z = \int e^{-U(x)} dx$$

↓

This satisfies the requirements of Boltzmann distribution  $\rightarrow$  lower energy states are more probable, consistent with laws of thermodynamics, statistical mechanics.  
Arises naturally from max. entropy.

# Euler - Maruyama Discretization (ULA)

## Unadjusted Langevin Algorithm

$$X_{k+1} = X_k - \gamma_{k+1} \nabla U(X_k) + \sqrt{2\gamma_{k+1}} Z_{k+1}$$

$X_k \rightarrow$  state at the  $k$ -th iteration

$\gamma_{k+1} \rightarrow$  step size at  $k+1$ th iteration.

$Z_{k+1} \rightarrow$  standard normal R.V.

$$Z_{k+1} \sim N(0, \mathbb{I})$$

$-\gamma_{k+1} \nabla U(X_k)$  - performs a gradient descent step towards regions of lower potential

$\sqrt{2\gamma_{k+1}} Z_{k+1} \rightarrow$  Adds Gaussian noise to ensure exploration.

step size  $\gamma_{k+1} \leftarrow$  controls trade off b/w

exploitation (small  $\gamma_{k+1}$  to move towards regions of lower potential) and exploration (to add noise to explore the space, i.e., large  $\gamma_{k+1}$ ).

In practice, a constant step size is often used.

Decreasing step size  $\gamma_{k+1} = \frac{c}{k+1} \rightarrow$  can improve

convergence but requires careful tuning.

$x_0 \rightarrow$  typically sampled from a simple dist.

initial state  $x_0 \sim \mathcal{N}(0, I)$  (Standard normal)

ULA assumes:  $\nabla U(x)$  is  $L$ -Lipschitz Continuous.

$$\| \nabla U(x) - \nabla U(y) \| \leq \underline{\underline{L}} \|x - y\| \quad \forall x, y$$

this assume the gradient doesn't change too rapidly which is necessary for the convergence of Euler-Maruyama scheme

$U(x) \rightarrow$  Smooth & differentiable.

Sampling :- ULA can be used to sample from target dist. of the form

$$p(x) = \frac{e^{-U(x)}}{Z} \rightarrow \text{normalizing constant which may/maynot be known.}$$

Algorithm 2:-

Initialize :-

- Start with an initial state  $x_0$

Iterate :-

For each iteration  $k$  :-

- Compute the gradient  $\nabla U(x_k)$

- Propose a new state :-

$$x_{k+1} = x_k - \gamma_{k+1} \nabla U(x_k) + \sqrt{2 \gamma_{k+1}} z_{k+1}$$

- Store  $x_{k+1}$  as sample.

- Burn-in - discard the first few samples to allow the chain to

converge

\* VLA doesn't need Metropolis-Hastings acceptance step, unlike MALA. Hence the samples may have some bias due to discretization error.

\* VLA - simpler & computationally cheaper than MALA.