# The NO-U-Turn Sample (**NUTS**)

Presented by Per Joachims on 19.12.2019

Research Seminar in Statistics: Bayesian computation: state of the art and recent developments

Humboldt-University of Berlin

# Introduction

▶ *Goal*: do bayesian inference based on posterior distribution

▶ *Problem*: posterior distribution often only known upto a normalizing constant:

$$p(\theta|x) \propto p(x|\theta)p(\theta)$$

▶ *Solution*: approximate posterior by **sampling** with Markov Chain Monte Carlo (MCMC)
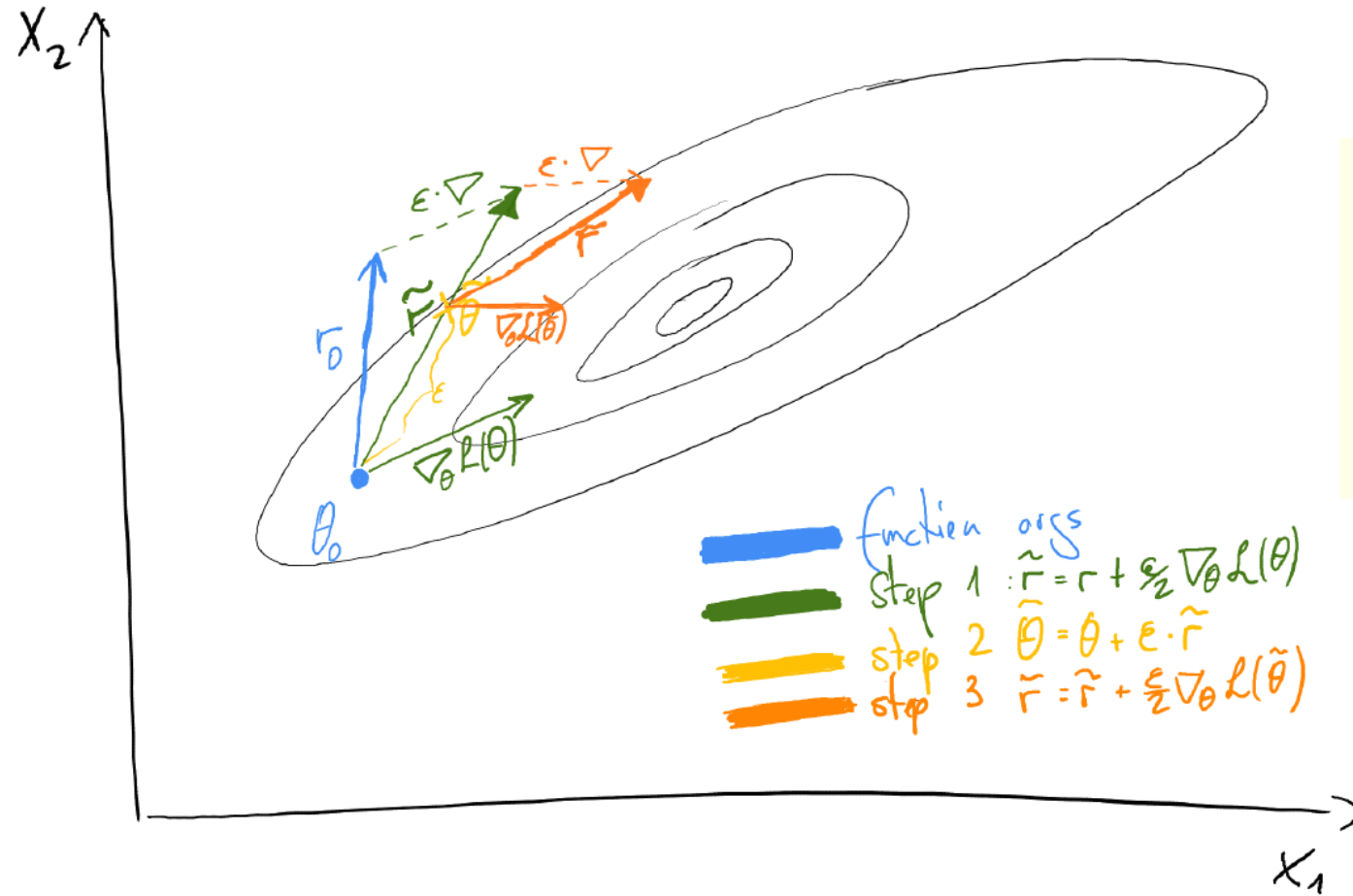
# NUTS and HMC

▶ The NO-U-Turn Sample (**NUTS**) by M.D. Hoffmann and A. Gelman (2011) is an extension of Hamilton Monte Carlo (**HMC**), which is a MCMC method

▶ HMC generates samples by setting up and simulating Hamilton dynamics

▶ With **r** = momentum vector, do repeatedly:

1. Sample r from multivariate normal

2. Evolve $\theta$, r by simulating L steps of the dynamics of the system

3. Accept or reject sample (similar to Metropolis)

Intro and HMC

$$\textbf{function } \text{Leapfrog}(\theta, r, \epsilon)$$
$$\text{Set } \tilde{r} \leftarrow r + (\epsilon/2)\nabla_\theta \mathcal{L}(\theta).$$
$$\text{Set } \tilde{\theta} \leftarrow \theta + \epsilon \tilde{r}.$$
$$\text{Set } \tilde{r} \leftarrow \tilde{r} + (\epsilon/2)\nabla_\theta \mathcal{L}(\tilde{\theta}).$$
$$\textbf{return } \tilde{\theta}, \tilde{r}.$$

function args
step 1 : $\tilde{r} = r + \frac{\epsilon}{2}\nabla_\theta \mathcal{L}(\theta)$
step 2 $\tilde{\theta} = \theta + \epsilon \cdot \tilde{r}$
step 3 $\tilde{r} = \tilde{r} + \frac{\epsilon}{2}\nabla_\theta \mathcal{L}(\tilde{\theta})$
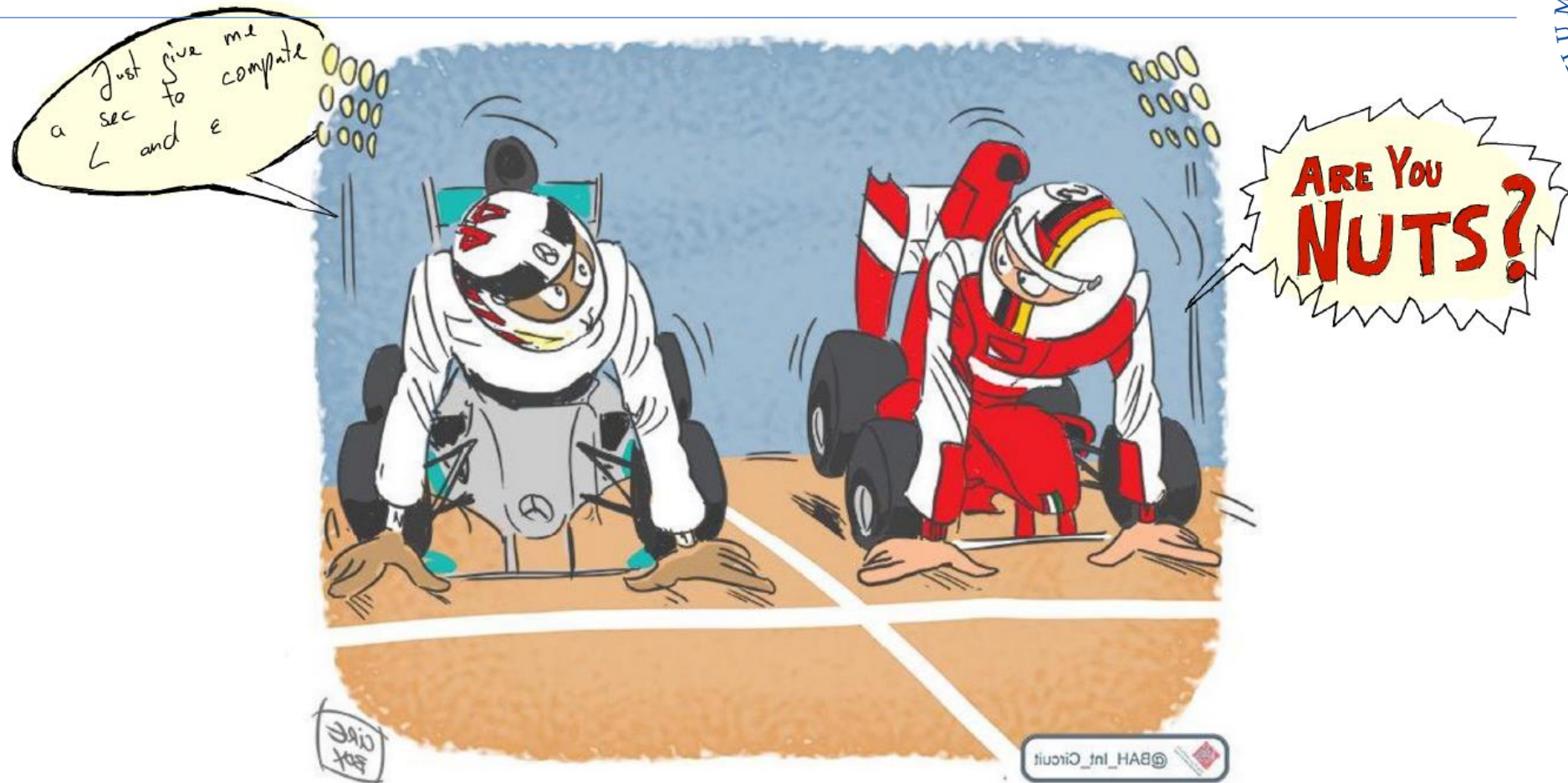
# HMC: Pro / Con

▶ *PRO:*

- **no** random walk behaviour which makes it (1) more **efficient** and (2) allows for better **dimension-scaling** compared to Metropolis Hastings and Gibbs samling

▶ *CON:*

- HMC needs to compute **gradients** -> not possible for discrecte variables

- Steps size ε and #steps L must be **tuned** well

Intro and HMC

Modified version of: https://twitter.com/sebvettelnews/status/589693881509253160

# NUTS in a Nutshell

- Eliminates the need to hand-tune L (and ε) and makes it **available** to more people

- Efficiency (NUTS) >= Efficiency (well tuned HMC)

- *Core Idea:* stop when the trajectory (path of gradient steps) starts to **turn back**

- *Why?* we do not want to "ruin" our progress of exploring the sample space made so far
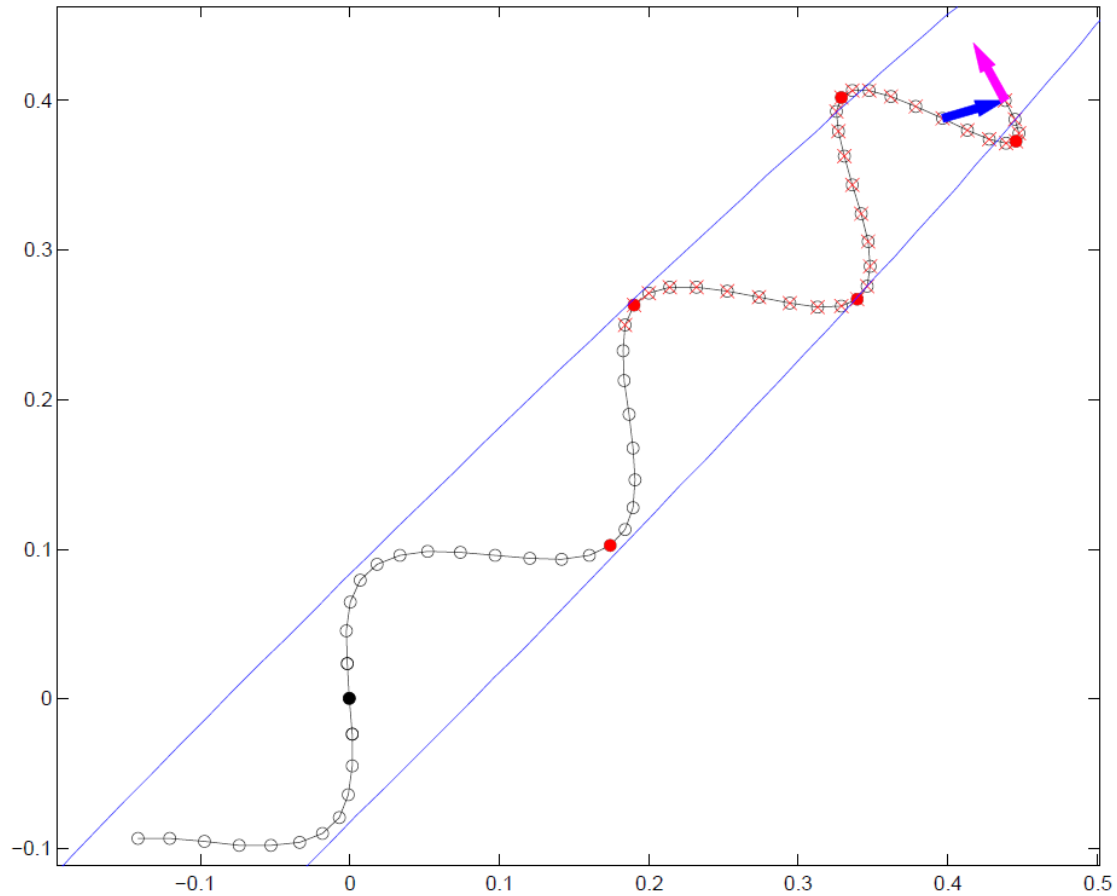
NUTS

# Definition of a U-Turn

▶ *Idea*: stop if the sampling iteration if the trajectory starts to turn back

▶ *Measure*: **Dot product!**

$$\frac{d}{dt}\frac{(\tilde{\theta}-\theta)\cdot(\tilde{\theta}-\theta)}{2} = (\tilde{\theta}-\theta)\cdot\frac{d}{dt}(\tilde{\theta}-\theta) = (\tilde{\theta}-\theta)\cdot\tilde{r}$$

▶ Sample backwards and forwards in time to fullfill reversibility condition

NUTS

# A Projectory Example



Legend

- ○    Positions on trajectory
- ●    Starting position
- ●    Excluded positions because of slicing
- ⊗    Excluded positions to satisfy detailed balance
- ➡ (magenta)    Momentum vector
- ➡ (blue)    Vector between states of subtree
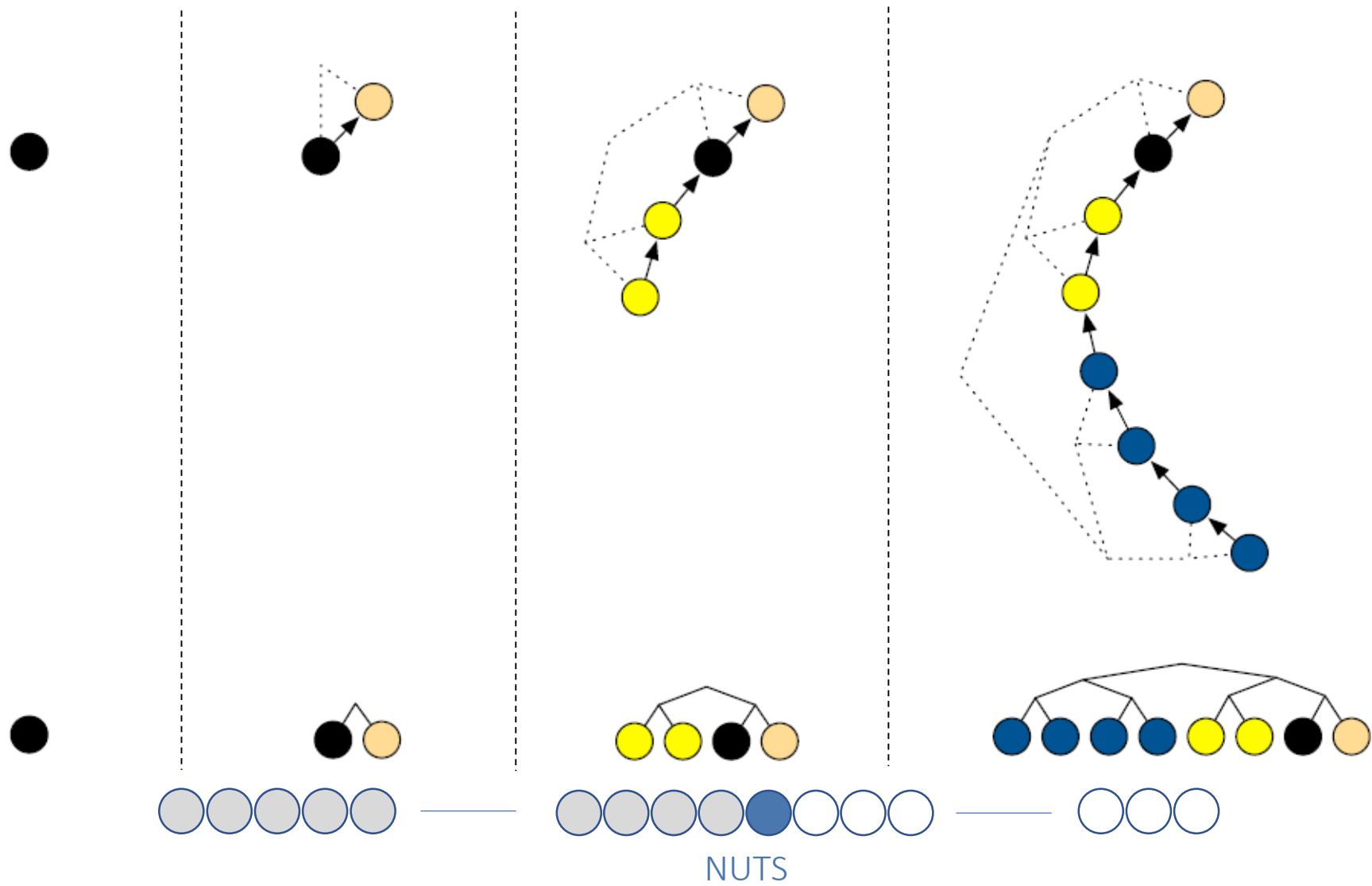
NUTS
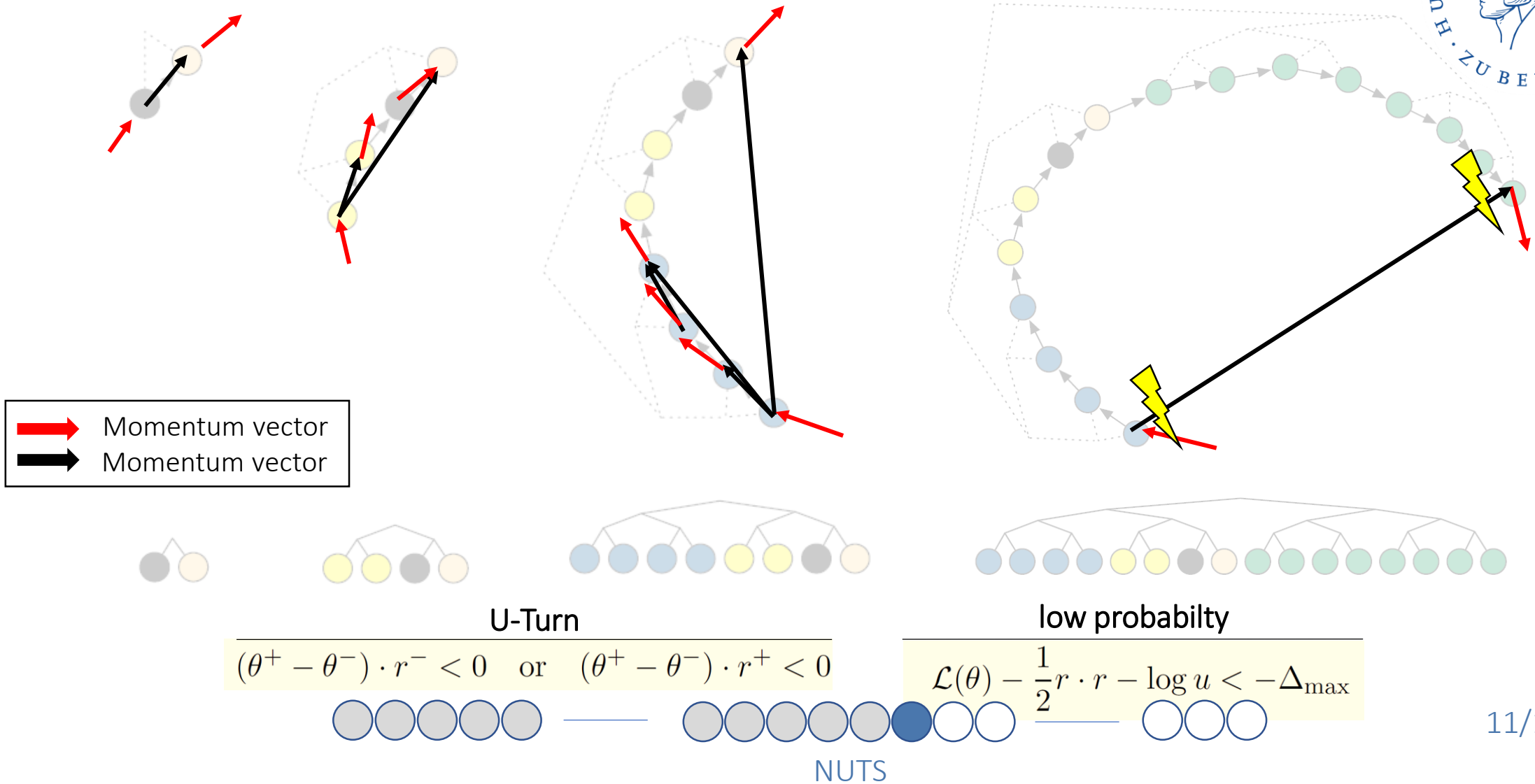
# The Naive NUTS (High Level)

▶ Sample **iteration**:

1. Sample momentum vector r

2. Sample slicing variable u

3. **Build trajectory:** trace out dynamics of θ, r forwards and backwards in time until stopping criteria (u-turn or/and probability checker) is matched

4. **Sample** uniformly from the points on trajectory

NUTS

NUTS

Momentum vector
Momentum vector

U-Turn

$$(\theta^+ - \theta^-) \cdot r^- < 0 \quad \text{or} \quad (\theta^+ - \theta^-) \cdot r^+ < 0$$

low probabilty

$$\mathcal{L}(\theta) - \frac{1}{2} r \cdot r - \log u < -\Delta_{\max}$$

NUTS

https://chi-feng.github.io/mcmc-demo/app.html
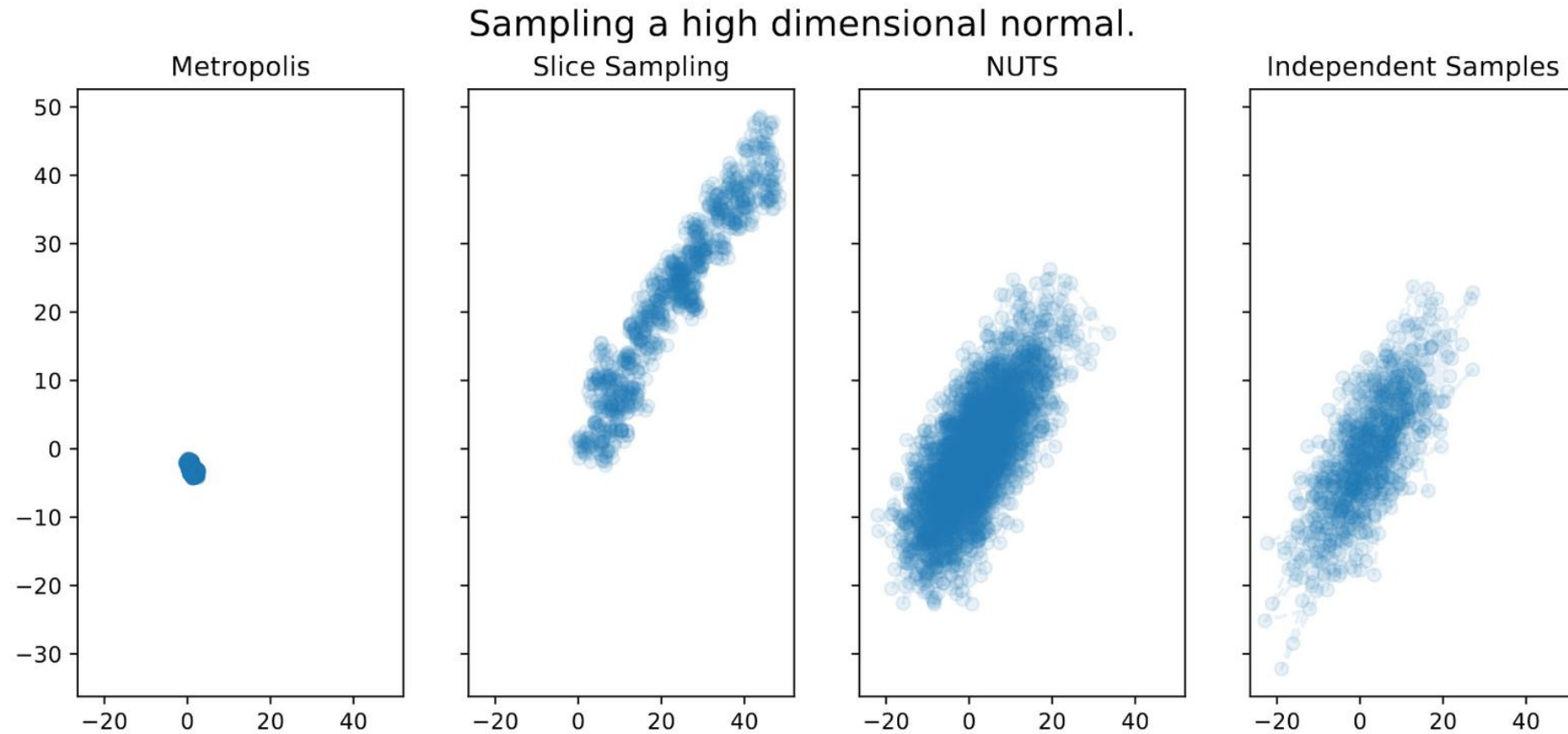
NUTS

▶ (1) Stop building the trajectory if **one** leaf of the tree matches low-probability-criterion -> save *computation*

▶ (2) Try to sample at the end of the trajectory -> more *efficient* exploration of the state space

▶ (3) **sample θ, r while building** the tree to reduce size of *memory* used: $O(j) < O(2^j)$

▶ (4) Set step size **ε** with dual averaging


ONE DOES NOT SIMPLY
SAMPLE UNIFORMLY AND WASTE COMPUTATION

NUTS
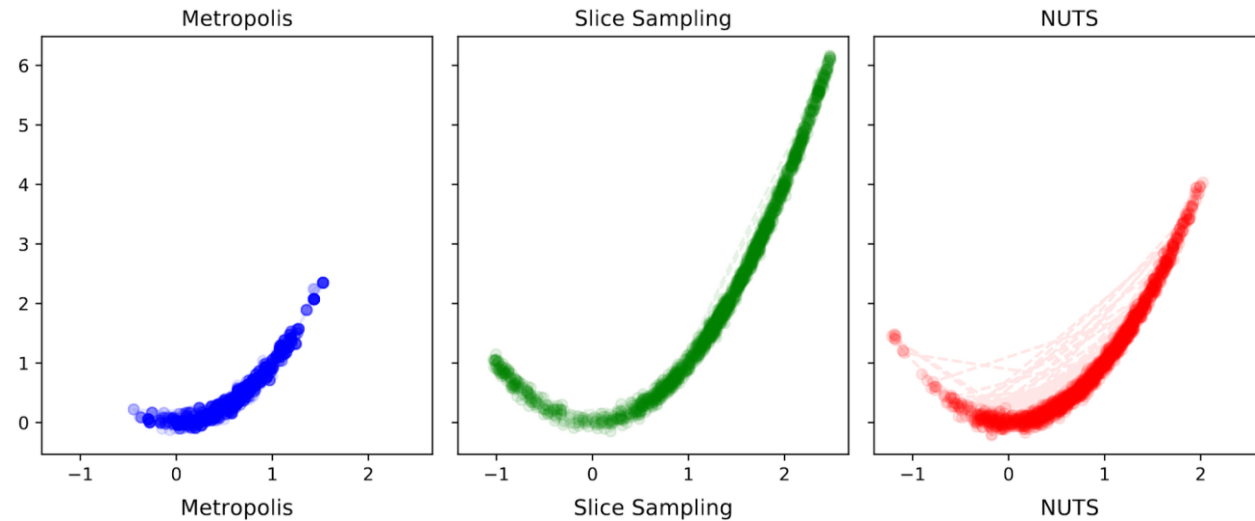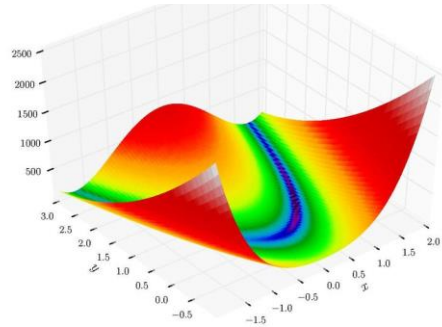
Sampling a high dimensional normal.

Ex. and Wrap-up

# NUTS: Perks and Limitations

**Rosenbrock's banana function**

$$f(x,y) = (a-x)^2 + b(y-x^2)^2$$

**Himmelblau's function**

$$f(x,y) = (x^2 + y - 11)^2 + (x + y^2 - 7)^2.$$

Ex. and Wrap-up

▶ *PRO*:

- no random walk behaviour which makes it (1) more **efficient** and (2) allows for better **dimension-scaling** compared to Metropolis Hastings and Gibbs samling

- no need for handtuning

▶ *CON:*

- as HMC, NUTS needs to compute **gradients** -> not possible for disrecte variables

➡ NUTS is default sampling method in STAN / pymc3

Ex. and Wrap-up

# THANKS !

See this + code https://github.com/pjoachims/rssbayes

# References / Further Readings

[1] Hoffman, M.D. and Gelman, A., 2014. The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research*, *15*(1), pp.1593-1623.

[2] Betancourt, M., 2017. A conceptual introduction to Hamiltonian Monte Carlo. *arXiv preprint arXiv:1701.02434*.

[3] Neal, R.M., 1994. An improved acceptance procedure for the hybrid Monte Carlo algorithm. *Journal of Computational Physics*, *111*(1), pp.194-203.


[4] https://chi-feng.github.io/mcmc-demo/

[5] https://docs.pymc.io/

[6] https://en.wikipedia.org/wiki/Test_functions_for_optimization

**Algorithm 1** Hamiltonian Monte Carlo

Given $\theta^0$, $\epsilon$, $L$, $\mathcal{L}$, $M$:

**for** $m = 1$ to $M$ **do**

    Sample $r^0 \sim \mathcal{N}(0, I)$.

    Set $\theta^m \leftarrow \theta^{m-1}, \tilde{\theta} \leftarrow \theta^{m-1}, \tilde{r} \leftarrow r^0$.

    **for** $i = 1$ to $L$ **do**

        Set $\tilde{\theta}, \tilde{r} \leftarrow \text{Leapfrog}(\tilde{\theta}, \tilde{r}, \epsilon)$.

    **end for**

    With probability $\alpha = \min\left\{1, \frac{\exp\{\mathcal{L}(\tilde{\theta}) - \frac{1}{2}\tilde{r}\cdot\tilde{r}\}}{\exp\{\mathcal{L}(\theta^{m-1}) - \frac{1}{2}r^0\cdot r^0\}}\right\}$, set $\theta^m \leftarrow \tilde{\theta}$, $r^m \leftarrow -\tilde{r}$.

**end for**

**function** Leapfrog$(\theta, r, \epsilon)$

Set $\tilde{r} \leftarrow r + (\epsilon/2)\nabla_\theta \mathcal{L}(\theta)$.

Set $\tilde{\theta} \leftarrow \theta + \epsilon\tilde{r}$.

Set $\tilde{r} \leftarrow \tilde{r} + (\epsilon/2)\nabla_\theta \mathcal{L}(\tilde{\theta})$.

**return** $\tilde{\theta}, \tilde{r}$.

**Algorithm 2** Naive No-U-Turn Sampler

Given $\theta^0$, $\epsilon$, $\mathcal{L}$, $M$:

for $m = 1$ to $M$ do

    Resample $r^0 \sim \mathcal{N}(0, I)$.

    Resample $u \sim \text{Uniform}([0, \exp\{\mathcal{L}(\theta^{m-1} - \frac{1}{2}r^0 \cdot r^0\}])$

    Initialize $\theta^- = \theta^{m-1}$, $\theta^+ = \theta^{m-1}$, $r^- = r^0$, $r^+ = r^0$, $j = 0$, $\mathcal{C} = \{(\theta^{m-1}, r^0)\}$, $s = 1$.

    while $s = 1$ do

        Choose a direction $v_j \sim \text{Uniform}(\{-1, 1\})$.

        if $v_j = -1$ then

            $\theta^-, r^-, -, -, \mathcal{C}', s' \leftarrow \text{BuildTree}(\theta^-, r^-, u, v_j, j, \epsilon)$.

        else

            $-, -, \theta^+, r^+, \mathcal{C}', s' \leftarrow \text{BuildTree}(\theta^+, r^+, u, v_j, j, \epsilon)$.

        end if

        if $s' = 1$ then

            $\mathcal{C} \leftarrow \mathcal{C} \cup \mathcal{C}'$.

        end if

        $s \leftarrow s' \mathbb{I}[(\theta^+ - \theta^-) \cdot r^- \geq 0] \mathbb{I}[(\theta^+ - \theta^-) \cdot r^+ \geq 0]$.

        $j \leftarrow j + 1$.

    end while

    Sample $\theta^m, r$ uniformly at random from $\mathcal{C}$.

end for

---

function $\text{BuildTree}(\theta, r, u, v, j, \epsilon)$

if $j = 0$ then

    *Base case—take one leapfrog step in the direction $v$.*

    $\theta', r' \leftarrow \text{Leapfrog}(\theta, r, v\epsilon)$.

    $\mathcal{C}' \leftarrow \begin{cases} \{(\theta', r')\} & \text{if } u \leq \exp\{\mathcal{L}(\theta') - \frac{1}{2}r' \cdot r'\} \\ \emptyset & \text{else} \end{cases}$

    $s' \leftarrow \mathbb{I}[\mathcal{L}(\theta') - \frac{1}{2}r' \cdot r' > \log u - \Delta_{\max}]$.

    return $\theta', r', \theta', r', \mathcal{C}', s'$.

else

    *Recursion—build the left and right subtrees.*

    $\theta^-, r^-, \theta^+, r^+, \mathcal{C}', s' \leftarrow \text{BuildTree}(\theta, r, u, v, j - 1, \epsilon)$.

    if $v = -1$ then

        $\theta^-, r^-, -, -, \mathcal{C}'', s'' \leftarrow \text{BuildTree}(\theta^-, r^-, u, v, j - 1, \epsilon)$.

    else

        $-, -, \theta^+, r^+, \mathcal{C}'', s'' \leftarrow \text{BuildTree}(\theta^+, r^+, u, v, j - 1, \epsilon)$.

    end if

    $s' \leftarrow s' s'' \mathbb{I}[(\theta^+ - \theta^-) \cdot r^- \geq 0] \mathbb{I}[(\theta^+ - \theta^-) \cdot r^+ \geq 0]$.

    $\mathcal{C}' \leftarrow \mathcal{C}' \cup \mathcal{C}''$.

    return $\theta^-, r^-, \theta^+, r^+, \mathcal{C}', s'$.

end if

**Algorithm 3** Efficient No-U-Turn Sampler

Given $\theta^0$, $\epsilon$, $\mathcal{L}$, $M$:

for $m = 1$ to $M$ do

    Resample $r^0 \sim \mathcal{N}(0, I)$.

    Resample $u \sim \text{Uniform}([0, \exp\{\mathcal{L}(\theta^{m-1} - \frac{1}{2}r^0 \cdot r^0\}])$

    Initialize $\theta^- = \theta^{m-1}$, $\theta^+ = \theta^{m-1}$, $r^- = r^0$, $r^+ = r^0$, $j = 0$, $\theta^m = \theta^{m-1}$, $n = 1$, $s = 1$.

    while $s = 1$ do

        Choose a direction $v_j \sim \text{Uniform}(\{-1, 1\})$.

        if $v_j = -1$ then

            $\theta^-, r^-, -, -, \theta', n', s' \leftarrow \text{BuildTree}(\theta^-, r^-, u, v_j, j, \epsilon)$.

        else

            $-, -, \theta^+, r^+, \theta', n', s' \leftarrow \text{BuildTree}(\theta^+, r^+, u, v_j, j, \epsilon)$.

        end if

        if $s' = 1$ then

            With probability $\min\{1, \frac{n'}{n}\}$, set $\theta^m \leftarrow \theta'$.

        end if

        $n \leftarrow n + n'$.

        $s \leftarrow s'\mathbb{I}[(\theta^+ - \theta^-) \cdot r^- \geq 0]\mathbb{I}[(\theta^+ - \theta^-) \cdot r^+ \geq 0]$.

        $j \leftarrow j + 1$.

    end while

end for

---

function BuildTree($\theta, r, u, v, j, \epsilon$)

if $j = 0$ then

    *Base case—take one leapfrog step in the direction $v$.*

    $\theta', r' \leftarrow \text{Leapfrog}(\theta, r, v\epsilon)$.

    $n' \leftarrow \mathbb{I}[u \leq \exp\{\mathcal{L}(\theta') - \frac{1}{2}r' \cdot r'\}]$.

    $s' \leftarrow \mathbb{I}[\mathcal{L}(\theta') - \frac{1}{2}r' \cdot r' > \log u - \Delta_{\max}]$

    return $\theta', r', \theta', r', \theta', n', s'$.

else

    *Recursion—implicitly build the left and right subtrees.*

    $\theta^-, r^-, \theta^+, r^+, \theta', n', s' \leftarrow \text{BuildTree}(\theta, r, u, v, j - 1, \epsilon)$.

    if $s' = 1$ then

        if $v = -1$ then

            $\theta^-, r^-, -, -, \theta'', n'', s'' \leftarrow \text{BuildTree}(\theta^-, r^-, u, v, j - 1, \epsilon)$.

        else

            $-, -, \theta^+, r^+, \theta'', n'', s'' \leftarrow \text{BuildTree}(\theta^+, r^+, u, v, j - 1, \epsilon)$.

        end if

        With probability $\frac{n''}{n'+n''}$, set $\theta' \leftarrow \theta''$.

        $s' \leftarrow s''\mathbb{I}[(\theta^+ - \theta^-) \cdot r^- \geq 0]\mathbb{I}[(\theta^+ - \theta^-) \cdot r^+ \geq 0]$

        $n' \leftarrow n' + n''$

    end if

    return $\theta^-, r^-, \theta^+, r^+, \theta', n', s'$.

end if

C.1: All elements of $\mathcal{C}$ must be chosen in a way that preserves volume. That is, any deterministic transformations of $\theta, r$ used to add a state $\theta', r'$ to $\mathcal{C}$ must have a Jacobian with unit determinant.

Why?: threat unnormalized prob. density of element as unconditional probaility mass
How?: leapfrog step are volume preserving

C.2: $p((\theta, r) \in \mathcal{C} | \theta, r, u, \epsilon) = 1.$

Why?: ensure current state is a valid sample
How?: satisfied if inital state is in proposal states

C.3: $p(u \leq \exp\{\mathcal{L}(\theta') - \frac{1}{2} r' \cdot r'\} | (\theta', r') \in \mathcal{C}) = 1.$

Why?: any state in C must be in the slice defined by u -> all states have equal and positve cond. prob. density
How?: satisfied due to the slicing variable

C.4: If $(\theta, r) \in \mathcal{C}$ and $(\theta', r') \in \mathcal{C}$ then for any $\mathcal{B}$, $p(\mathcal{B}, \mathcal{C} | \theta, r, u, \epsilon) = p(\mathcal{B}, \mathcal{C} | \theta', r', u, \epsilon).$

Why?: B and C must have equal prob. to be selected regardless of the state
How?: exclude from C any state that could not have generated B

$$T(w'|w, \mathcal{C}) = \begin{cases} \dfrac{\mathbb{I}[w' \in \mathcal{C}^{\mathrm{new}}]}{|\mathcal{C}^{\mathrm{new}}|} & \text{if } |\mathcal{C}^{\mathrm{new}}| > |\mathcal{C}^{\mathrm{old}}|, \\[2ex] \dfrac{|\mathcal{C}^{\mathrm{new}}|}{|\mathcal{C}^{\mathrm{old}}|} \dfrac{\mathbb{I}[w' \in \mathcal{C}^{\mathrm{new}}]}{|\mathcal{C}^{\mathrm{new}}|} + \left(1 - \dfrac{|\mathcal{C}^{\mathrm{new}}|}{|\mathcal{C}^{\mathrm{old}}|}\right) \mathbb{I}[w' = w] & \text{if } |\mathcal{C}^{\mathrm{new}}| \le |\mathcal{C}^{\mathrm{old}}| \end{cases}$$

**Algorithm 2** Naive No-U-Turn Sampler

Given $\theta^0$, $\epsilon$, $\mathcal{L}$, $M$:

for $m = 1$ to $M$ do

    Resample $r^0 \sim \mathcal{N}(0, I)$.

    Resample $u \sim \mathrm{Uniform}([0, \exp\{\mathcal{L}(\theta^{m-1} - \frac{1}{2}r^0 \cdot r^0\}])$

    Initialize $\theta^- = \theta^{m-1}$, $\theta^+ = \theta^{m-1}$, $r^- = r^0$, $r^+ = r^0$, $j = 0$, $\mathcal{C} = \{(\theta^{m-1}, r^0)\}$, $s = 1$.

    while $s = 1$ do

        Choose a direction $v_j \sim \mathrm{Uniform}(\{-1, 1\})$.

        if $v_j = -1$ then

            $\theta^-, r^-, -, -, \mathcal{C}', s' \leftarrow \mathrm{BuildTree}(\theta^-, r^-, u, v_j, j, \epsilon)$.

        else

            $-, -, \theta^+, r^+, \mathcal{C}', s' \leftarrow \mathrm{BuildTree}(\theta^+, r^+, u, v_j, j, \epsilon)$.

        end if

        if $s' = 1$ then

            $\mathcal{C} \leftarrow \mathcal{C} \cup \mathcal{C}'$.

        end if

        $s \leftarrow s' \mathbb{I}[(\theta^+ - \theta^-) \cdot r^- \geq 0] \mathbb{I}[(\theta^+ - \theta^-) \cdot r^+ \geq 0]$.

        $j \leftarrow j + 1$.

    end while

    Sample $\theta^m, r$ uniformly at random from $\mathcal{C}$.

end for

**Algorithm 3** Efficient No-U-Turn Sampler

Given $\theta^0$, $\epsilon$, $\mathcal{L}$, $M$:

for $m = 1$ to $M$ do

    Resample $r^0 \sim \mathcal{N}(0, I)$.

    Resample $u \sim \mathrm{Uniform}([0, \exp\{\mathcal{L}(\theta^{m-1} - \frac{1}{2}r^0 \cdot r^0\}])$

    Initialize $\theta^- = \theta^{m-1}$, $\theta^+ = \theta^{m-1}$, $r^- = r^0$, $r^+ = r^0$, $j = 0$, $\theta^m = \theta^{m-1}$, $n = 1$, $s = 1$.

    while $s = 1$ do

        Choose a direction $v_j \sim \mathrm{Uniform}(\{-1, 1\})$.

        if $v_j = -1$ then

            $\theta^-, r^-, -, -, \theta', n', s' \leftarrow \mathrm{BuildTree}(\theta^-, r^-, u, v_j, j, \epsilon)$.

        else

            $-, -, \theta^+, r^+, \theta', n', s' \leftarrow \mathrm{BuildTree}(\theta^+, r^+, u, v_j, j, \epsilon)$.

        end if

        if $s' = 1$ then

            With probability $\min\{1, \frac{n'}{n}\}$, set $\theta^m \leftarrow \theta'$.

        end if

        $n \leftarrow n + n'$.

        $s \leftarrow s' \mathbb{I}[(\theta^+ - \theta^-) \cdot r^- \geq 0] \mathbb{I}[(\theta^+ - \theta^-) \cdot r^+ \geq 0]$.

        $j \leftarrow j + 1$.

    end while

end for

## Naive

$\textbf{function BuildTree}(\theta, r, u, v, j, \epsilon)$
$\textbf{if } j = 0 \textbf{ then}$
    *Base case—take one leapfrog step in the direction v.*
    $\theta', r' \leftarrow \text{Leapfrog}(\theta, r, v\epsilon).$
    $C' \leftarrow \begin{cases} \{(\theta', r')\} & \text{if } u \leq \exp\{\mathcal{L}(\theta') - \frac{1}{2}r' \cdot r'\} \\ \emptyset & \text{else} \end{cases}$
    $s' \leftarrow \mathbb{I}[\mathcal{L}(\theta') - \frac{1}{2}r' \cdot r' > \log u - \Delta_{\max}].$
    $\textbf{return } \theta', r', \theta', r', C', s'.$
$\textbf{else}$
    *Recursion—build the left and right subtrees.*
    $\theta^-, r^-, \theta^+, r^+, C', s' \leftarrow \text{BuildTree}(\theta, r, u, v, j - 1, \epsilon).$
    $\textbf{if } v = -1 \textbf{ then}$
        $\theta^-, r^-, -, -, C'', s'' \leftarrow \text{BuildTree}(\theta^-, r^-, u, v, j - 1, \epsilon).$
    $\textbf{else}$
        $-, -, \theta^+, r^+, C'', s'' \leftarrow \text{BuildTree}(\theta^+, r^+, u, v, j - 1, \epsilon).$
    $\textbf{end if}$
    $s' \leftarrow s's''\mathbb{I}[(\theta^+ - \theta^-) \cdot r^- \geq 0]\mathbb{I}[(\theta^+ - \theta^-) \cdot r^+ \geq 0].$
    $C' \leftarrow C' \cup C''.$
    $\textbf{return } \theta^-, r^-, \theta^+, r^+, C', s'.$
$\textbf{end if}$

## Efficient

$\textbf{function BuildTree}(\theta, r, u, v, j, \epsilon)$
$\textbf{if } j = 0 \textbf{ then}$
    *Base case—take one leapfrog step in the direction v.*
    $\theta', r' \leftarrow \text{Leapfrog}(\theta, r, v\epsilon).$
    $n' \leftarrow \mathbb{I}[u \leq \exp\{\mathcal{L}(\theta') - \frac{1}{2}r' \cdot r'\}].$
    $s' \leftarrow \mathbb{I}[\mathcal{L}(\theta') - \frac{1}{2}r' \cdot r' > \log u - \Delta_{\max}]$
    $\textbf{return } \theta', r', \theta', r', \theta', n', s'.$
$\textbf{else}$
    *Recursion—implicitly build the left and right subtrees.*
    $\theta^-, r^-, \theta^+, r^+, \theta', n', s' \leftarrow \text{BuildTree}(\theta, r, u, v, j - 1, \epsilon).$
    $\textbf{if } s' = 1 \textbf{ then}$
        $\textbf{if } v = -1 \textbf{ then}$
            $\theta^-, r^-, -, -, \theta'', n'', s'' \leftarrow \text{BuildTree}(\theta^-, r^-, u, v, j - 1, \epsilon).$
        $\textbf{else}$
            $-, -, \theta^+, r^+, \theta'', n'', s'' \leftarrow \text{BuildTree}(\theta^+, r^+, u, v, j - 1, \epsilon).$
        $\textbf{end if}$
        With probability $\frac{n''}{n'+n''}$, set $\theta' \leftarrow \theta''.$
        $s' \leftarrow s''\mathbb{I}[(\theta^+ - \theta^-) \cdot r^- \geq 0]\mathbb{I}[(\theta^+ - \theta^-) \cdot r^+ \geq 0]$
        $n' \leftarrow n' + n''$
    $\textbf{end if}$
    $\textbf{return } \theta^-, r^-, \theta^+, r^+, \theta', n', s'.$
$\textbf{end if}$