

A RV is a function that maps an event in sample space
to a number

$$X: \Omega \mapsto \mathbb{R}$$

2 coin tosses $\rightarrow X$ is # heads

HH	HT
TH	TT

$$X = \begin{cases} HH \rightarrow 2 \\ HT \rightarrow 1 \\ TH \rightarrow 1 \\ TT \rightarrow 0 \end{cases}$$

Family of distributions

$$N(\mu, \sigma^2)$$

Every combination of $\underline{\theta} = (\mu, \sigma)$ defines a member of the normal family $\rightarrow \theta = (\mu, \sigma)$ are parameters of the distribution

Types of distribution

Discrete $X \in \mathbb{Z}$



Continuous $X \in \mathbb{R}$



Distributions



E.g. Suppose IQ is normally distributed with

$$\mu = 100 \quad \sigma = 15$$

- ① Fraction $\equiv \text{IQ} \leq 120$, $80 \leq \text{IQ} \leq 110$ (cdf)
- ② IQ to be in top 5% of population (quantile)
- ③ Mode of IQ (pdf)
- ④ Sample 100 IQ values in top 25 percent (rvs)

$$\begin{cases} \text{norm.cdf} \equiv \text{pnorm} & \text{norm.pdf} \equiv \text{dnorm} \\ \text{norm.pdf} \equiv \text{dnorm} & \text{norm.rvs} \equiv \text{rnorm} \end{cases}$$

scipy.stats

Likelihood

$$\begin{array}{ll} \text{Probability} & x \mapsto f(x|\theta) \\ \text{Likelihood} & \theta \mapsto f(x|\theta) \\ L(\theta|x) = f(x|\theta) \end{array} \quad \left. \right\} \text{Equivalent}$$

For independent events, likelihood of n events is
the product of the individual likelihoods

$$\log \prod \downarrow \sum (\log \text{likelihood})$$
$$l(\theta|x)$$

$$\downarrow$$
$$\text{Funct.} \in \text{some min, max as } l(\theta|x)$$

For optimization,
often work in negative log likelihood

because we want a Minimization problem

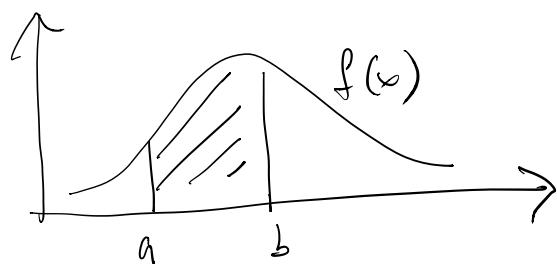
Maximum likelihood estimation (MLE)

Find θ where $L(\theta|x)$ is max }
≡ Find θ where $-L(\theta|x)$ is min } θ_{MLE}

$$\boxed{\theta_{MLE} = \arg \max L(\theta|x)}$$

→ Notebook 8A

Probability \Rightarrow Integration



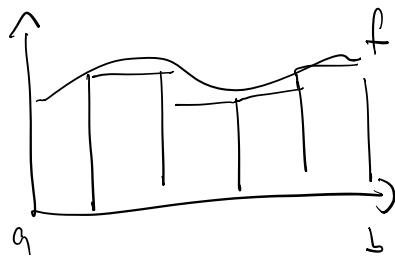
$$E[X] = \int x p(x) dx$$

$$E[f(x)] = \int f(x) p(x) dx$$

$$P(a \leq x \leq b) = \int_a^b f(x) dx$$

Finding probabilities
involves integration

Numerical Quadrature



Some form of
divide \Rightarrow conquer summation

Trapezoidal, Simpson's etc.

Scales poorly to \uparrow dimensions

if n is the # grid points $\rightarrow n^d$

$$n=10 \quad d=1 \quad 10$$

$$d=2 \quad 100$$

$$d=3 \quad 1,000$$

$$d=4 \quad 10,000$$

$$d=100 \quad \dots \quad]$$

Monte Carlo integration

$$E[g(x)] = \int_x^p g(x) p(x) dx$$

Lotus

$$\bar{g} = \frac{1}{n} \sum_{i=1}^n g(x_i) \quad x_i \sim p$$

Convergence is $O(n^{1/2})$ independent of dimensionality

Some intuition

$$I = \int f(x) dx$$

$$E[g(x)] = \int g(x) p(x) dx$$

Choose $g(x) = \frac{f(x)}{p(x)}$

$$E[g(x)] = \int \frac{f(x)}{p(x)} p(x) dx$$

$$= \int f(x) dx$$

$$= I$$

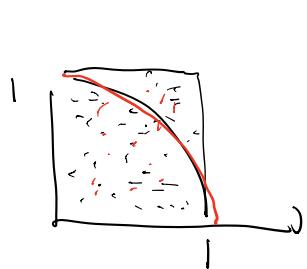
LLN

$$I \approx \bar{g} = \frac{1}{n} \sum_{i=1}^n g(x_i)$$

Estimating π .

$$f(x) = \begin{cases} 1 & \text{if } x^2 + y^2 \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$$I = \int_0^1 \int_0^1 f(x) dx = \frac{\pi}{4}$$



$$\hat{\pi} = 4 \cdot \frac{1}{N} \sum_{i=1}^N f(x_i)$$
$$x_i \sim U(0, 1)$$

Boyles theorem

Bayes theorem

$$p(\theta|x) = \frac{p(x|\theta) p(\theta)}{p(x)}$$

Parameter is RV
(c.f. frequentist)

posterior likelihood prior

$$= \frac{p(x|\theta) p(\theta)}{\int p(x|\theta) p(\theta) d\theta}$$

marginal
(normalizing constant)

often high-dimensional
? intractable

MAP

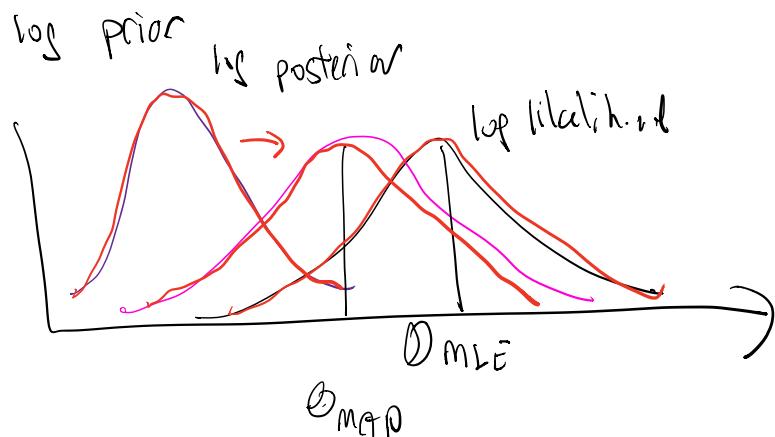
$$p(\theta|x) = \frac{p(x|\theta) p(\theta)}{\int p(x|\theta) p(\theta) d\theta}$$

Normalizing constant

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

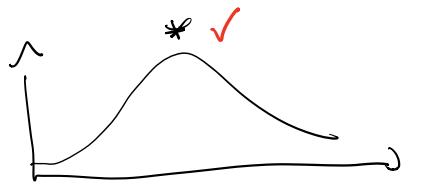
We can optimize \rightarrow just like maximum likelihood
 \simeq additional term for prior

\Rightarrow MAP estimate

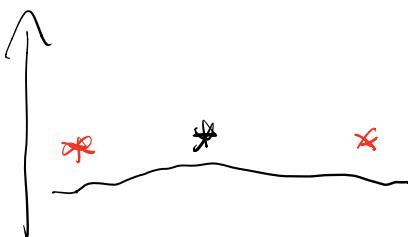


As $\eta \rightarrow \infty$ $\text{MAP} \rightarrow \text{MLE}$

Limitations of MAP, MLE

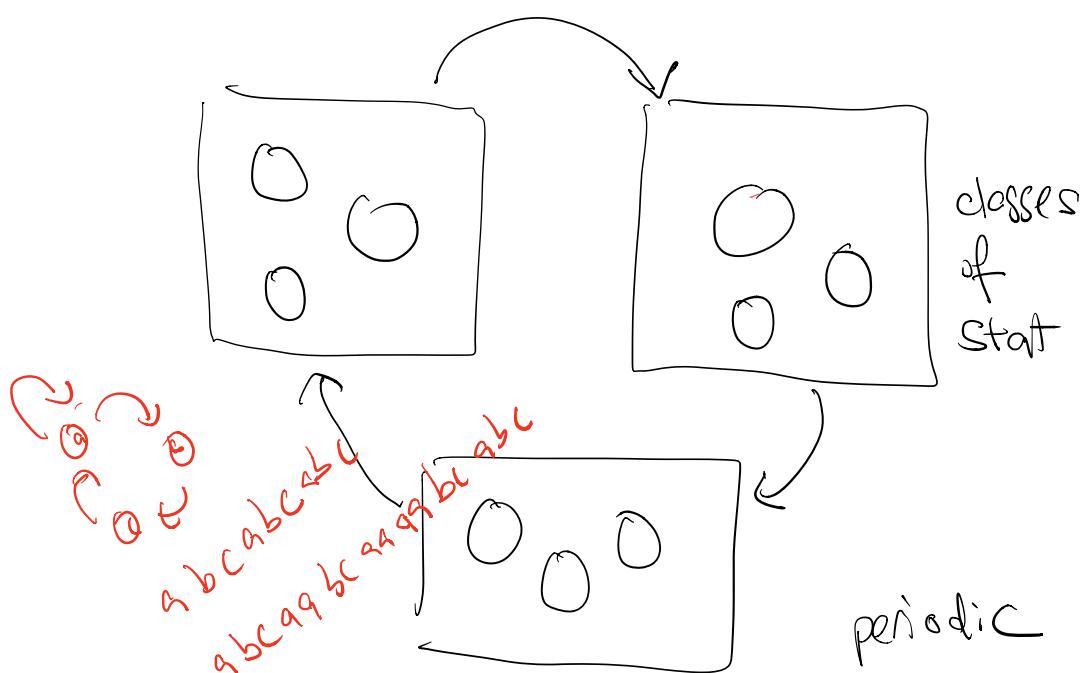
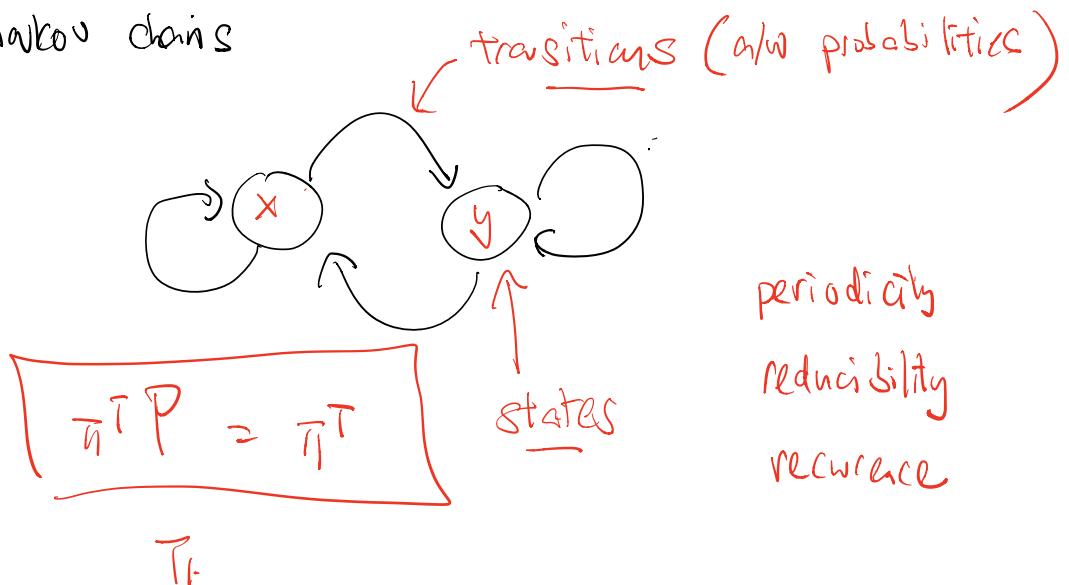


Same point estimate,
very different
posterior

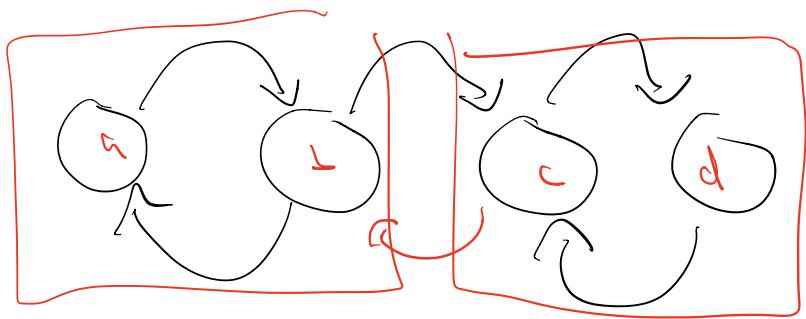


MCMC \rightarrow Markov Chain Monte Carlo

Markov chains



aperiodic / if not periodic



i, j communicate $\Leftrightarrow i \leftrightarrow j$

Put states into communication classes

reducible $\Rightarrow \geq 2$ communication classes

irreducible $\Rightarrow 1$ communication class

recurrent if $P(i \rightarrow i) > 0$

For continuous Markov chains

positive recurrent if expected

return time $< \infty$

Then as $t \rightarrow \infty$

i will be visited an ∞ # times

If a Markov chain is aperiodic, irreducible

↳ positive recurrent

⇒ unique steady state π

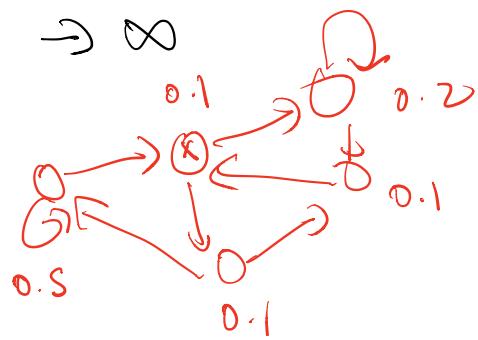
$$\pi^T P = \pi^T$$

Intuition

$\pi_i \approx$ Average # of times in state i

as $t \rightarrow \infty$

Ergodic



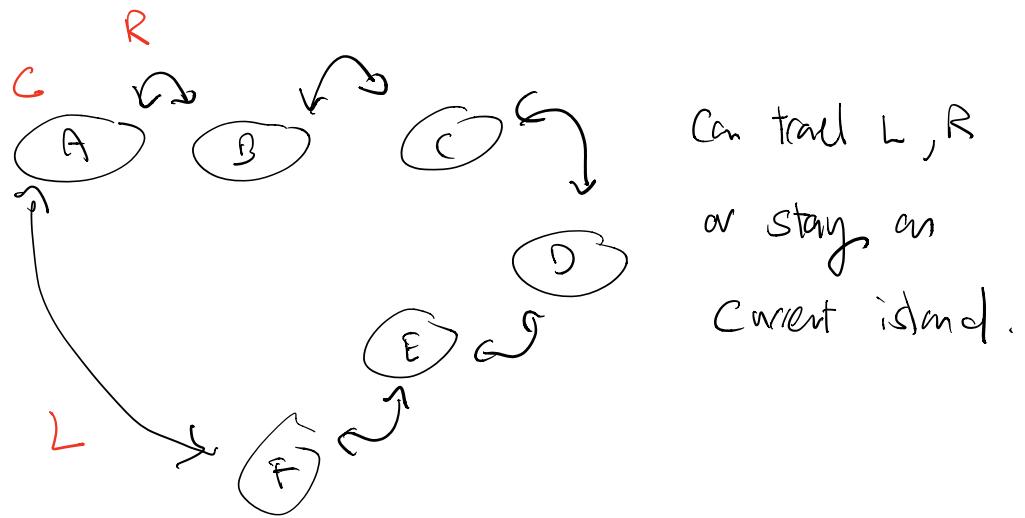
⑤ Why?

MCMC

①

steady state probability = posterior probability

Kruschke's island cosine example



Rules to move between islands

① Toss a coin

H → Ask for population (R) n_R

T → Ask for population (L) n_L

② Suppose you get H

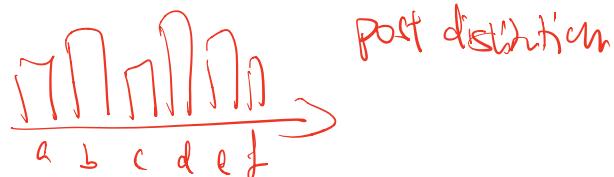
if $\underline{n_R} > \underline{n_C}$ → go (R)

else generate $p \in U(0, 1)$

If $p < \frac{n_R}{n_C}$ → go (R)

else → stay.

③ Repeat



④ Distribution of population in islands

& frequency of visits (over long time)

MH \rightarrow Metropolis Hastings symmetric proposal

① Toss a coin \leftarrow distribution

H \rightarrow Ask for population (R) n_R

T \rightarrow Ask for population (L) n_L

② Suppose you get H

if $n_R > n_C \rightarrow g_0(R)$ standard

else generate $p \in U(0, 1)$ \leftarrow uniform

If $p < \frac{n_R}{n_C} \rightarrow g_0(R)$

else stay

③ Repeat

+ weight
fnct. on

MH

(symmetric)

proposal distribution

θ (start)

$$\theta_p = \theta + \Delta\theta$$

$$\Delta\theta \sim N(\theta, \sigma)$$

$$\alpha = \frac{f(\theta_p | x)}{f(\theta | x)} \quad \begin{cases} \text{acceptance probability} \\ \text{Ratio of posterior} \\ \Rightarrow \text{marginal cancels out!} \end{cases}$$

$$\alpha = \frac{f(x|\theta_p) p(\theta_p)}{f(x|\theta) p(\theta)}$$

if $\alpha \geq 1$ $\theta \leftarrow$

if $\alpha < 1$ $\theta \leftarrow \theta_p$ = probability α

$\theta \leftarrow \theta$ = probability $1-\alpha$

$$\tilde{\theta} = \frac{1}{N} \sum_{i=1}^{k+N} \theta_i \quad \begin{array}{l} \rightarrow [\theta_0, \theta_1, \theta_2, \theta_2, \theta_2, \theta_3, \theta_4, \\ \theta_4] \\ \downarrow \\ \text{Burn-in} \end{array}$$

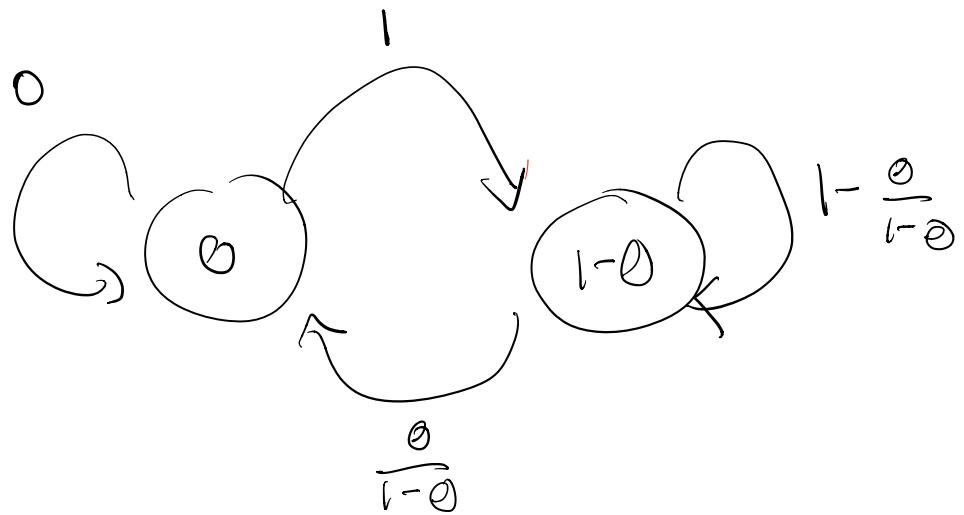
Why does MCMC work?

- ① There is a unique steady state
- ② By design, the steady state is the posterior distribution

System has just 2 states

Position probabilities

$$[\theta, 1-\theta] \quad \theta < 0.5$$



Solve stationary distribution

$$\begin{bmatrix} 0 & \theta/1-\theta \\ 1 & 1-\theta/1-\theta \end{bmatrix} \pi = \pi$$

$$\Rightarrow \pi = \begin{bmatrix} \theta \\ 1-\theta \end{bmatrix} \leftarrow$$

Want $\pi(x)$ to be posterior $P(x)$

$\boxed{0 \leftrightarrow 0}$

$$p(x)T(x \rightarrow y) = p(y)T(y \rightarrow x)$$

detailed balance \Rightarrow steady state

$$\frac{T(x \rightarrow y)}{T(y \rightarrow x)} = \frac{p(y)}{p(x)}$$

$$\frac{g(x \rightarrow y)A(x \rightarrow y)}{g(y \rightarrow x)A(y \rightarrow x)} = \frac{p(y)}{p(x)}$$

$$\boxed{\frac{A(x \rightarrow y)}{A(y \rightarrow x)} = \frac{p(y)g(y \rightarrow x)}{p(x)g(x \rightarrow y)}}$$

Choose $A(x \rightarrow y) \rightarrow \min\left(1, \frac{p(y)g(y \rightarrow x)}{p(x)g(x \rightarrow y)}\right)$

\Rightarrow

$$\boxed{\frac{A(x \rightarrow y)}{A(y \rightarrow x)} = \frac{p(y)g(y \rightarrow x)}{p(x)g(x \rightarrow y)}} \quad \checkmark$$

Metropolis
Metropolis-Hastings

Gibbs

$$\underline{\theta} = (\underline{\theta}_1, \underline{\theta}_2, \dots, \underline{\theta}_k) \quad p(\underline{\theta}_1 | \underline{\theta}_2, \underline{\theta}_3, x)$$

Suppose we can draw from full conditional distributions

$$\text{Post } \theta_i \leftarrow p(\theta_i | \underline{\theta}_{-i}, x) \quad \begin{matrix} \text{Need to be able to} \\ \text{find this} \end{matrix}$$

We cycle through each parameter \rightarrow at each step,

the proposal is the posterior (for one parameter)

\Rightarrow Always accept ✓

$$\frac{p(y) g(y \rightarrow x)}{p(x) g(x \rightarrow y)} = 1$$

① $p(x_{-i}) = p(y_{-i})$ Apart from parameter i
 old } proposed state
 are identical

② $p(x_i | x_{-i}) p(x_{-i}) = p(x_i, x_{-i}) = p(x)$ By definition

So

$$\frac{p(y) g(y \rightarrow x)}{p(x) g(x \rightarrow y)} = \frac{p(y_i | y_{-i}) p(y_{-i}) p(x_i | x_{-i})}{p(x_i | x_{-i}) p(x_{-i}) p(y_i | y_{-i})}$$

Gibbs



$$\Delta\theta \sim N(0, \varsigma)$$

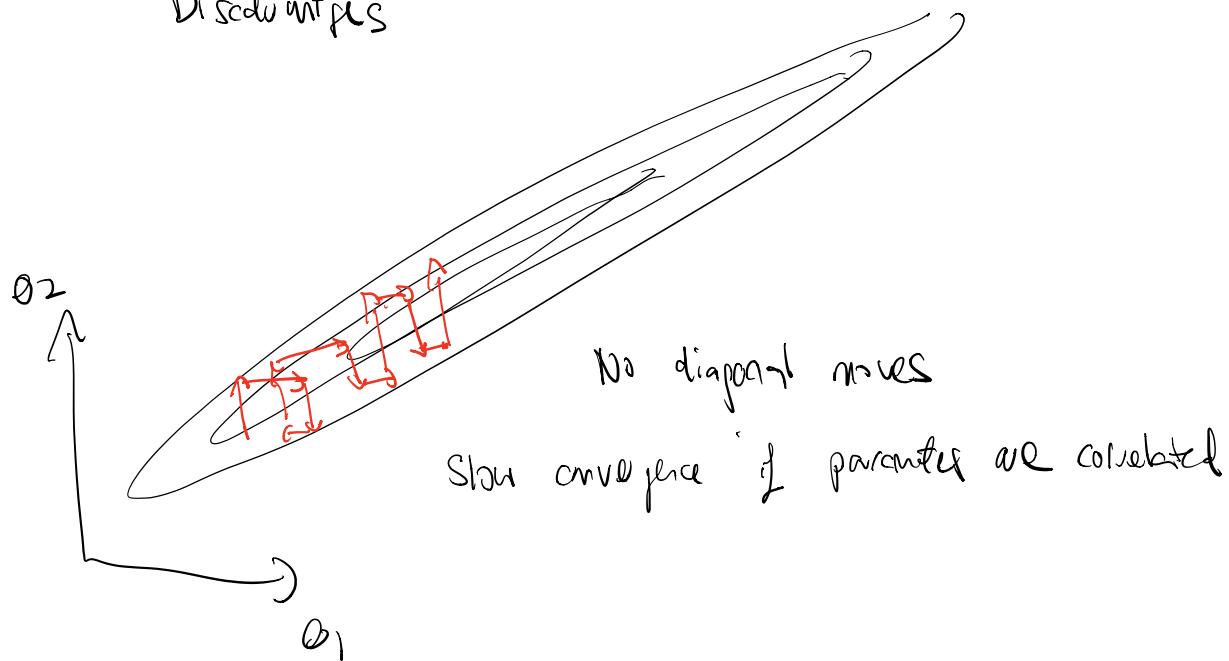
$$\varsigma \ll 1 \quad \bar{\nabla} E X$$

Advantage

① Always accept

② No tuning of proposal $\epsilon \gg 1 \quad \downarrow$

Disadvantages



Hamiltonian MC

From physics \rightarrow energy conservation, path independence

E.g. Frictionless well

$$KE + PE = \text{constant}$$

Particle $\bullet \rightarrow$ position \xrightarrow{x} velocity \xrightarrow{v} } state
 w/ unit mass

$$H(x, v) = K(v) + U(x)$$

Probability of state

$$p(x, v) \propto e^{-H(x, v)}$$

$$= e^{-U(x) - K(v)}$$

$$= e^{-U(x)} e^{-K(v)}$$

$$\propto p(x) p(v)$$

$$\Rightarrow p(x, v) = p(x) p(v)$$

x, v independent

Context with

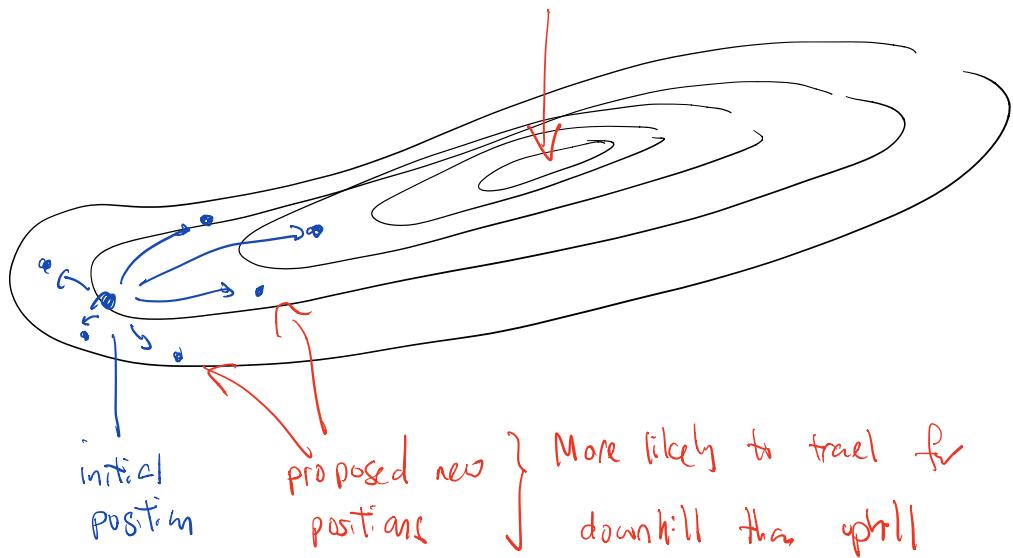
$$\left. \begin{aligned} p(x, v) &= p(v|x) p(x) \\ &= p(x|v) p(v) \end{aligned} \right\}$$

What happens to v does not affect x

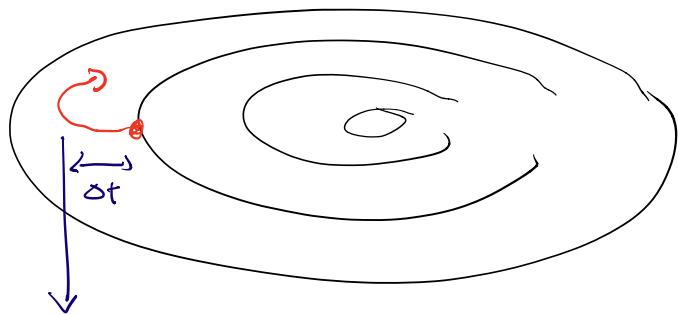
Proposal : Choose random v for particle, interpret
found in time, then use x^* as purpose

Acceptance ratio

$$\frac{e^{-U(x^*) - k(v^*)}}{e^{-U(x) - k(x)}} = e^{\underbrace{U(x) - U(x^*)}_{\text{minimum}} + \underbrace{k(v) - k(v^*)}_{\text{no numerical error.}}}$$



How long do we let particle travel?



Stop when there is a U-turn



No U-turn sampler (NUTS)