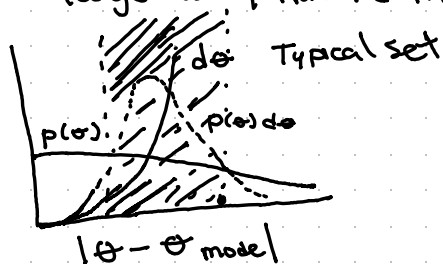# Why samples?

Maximum likelihood is not enough for many applications
or maximum a posteriori
Especially in high dimensions

The probability contained in a given region of param.
space $N$ is
$$\int_V p(\theta)\, d\theta$$

prob. density × volume

$p(\theta)$ peaks at the mode, but $d\theta$ is much
larger away from the mode in high-dimensions



$d\theta$ Typical set

$|\theta - \theta_{\text{mode}}|$

Samples allow us to compute expectation values

consider $X = f(\theta)$

$$\mathbb{E}[X] = \int_V f(\theta)\, p(\theta)\, d\theta$$

$$\sigma^2[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

If we have samples $\theta \sim p(\theta)$, then we can
estimate
$$\int_V f(\theta)\, p(\theta)\, d\theta \approx \langle f(\theta) \rangle_{\theta \sim p(\theta)} = \frac{1}{N}\sum_{i=1}^{N} f(\theta_i)$$

This __Monte Carlo__ estimate of the integral is unbiased:

$$\mathbb{E}[\langle f \rangle] = \frac{1}{N}\sum_{i=1}^{N} \mathbb{E}[f(\theta_i)] = \frac{1}{N}\sum_{i=1}^{N}\int_V f(\theta)\, p(\theta)\, d\theta$$

$$= \int_V f(\theta)\, p(\theta)\, d\theta \quad \checkmark$$

and converges with $\sqrt{N}$:

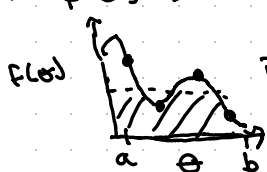$$\sigma^2[\langle f \rangle] = \sigma^2\left[\frac{1}{N}\sum_{i=1}^{N} f(\theta_i)\right]$$

$$= \frac{1}{N^2}\sum_{i=1}^{N}\sigma^2[f(\theta_i)] \qquad \text{if samples are independent, variance of sum is sum of variance}$$

$$= \frac{1}{N}\sigma^2[f(\theta)]$$

if samples are _not_ independent, replace $N$
with "effective sample size"

Ex: if $p(\theta)$ is uniform,



$f(\theta)$

$$\frac{1}{b-a}\int_a^b f(\theta)\, d\theta \approx \langle f(\theta) \rangle \cdot \frac{(b-a)}{b-a}$$

How can we get $\{\vec{\theta}_i\}$ samples from $p(\vec{\theta})$?
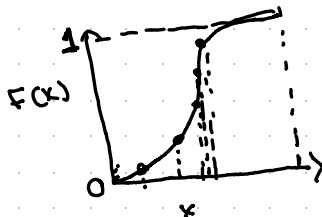
- inverse transform sampling from target with pdf $p(x)$

In 1-d, CDF $F(x) = \int_0^x p(x') \, dx'$

$\underbrace{\phantom{p(x')}}_{pdf}$

$0 < F(x) < 1$

draw $y_i = F(x_i)$ from uniform distribution

$y \sim U(0,1)$

$F^{-1}(y_i) = x_i$    $F(x)$



- rejection sampling from target with pdf $p(x)$

draw samples from some easy-to-sample distribution

$x_i \sim g(x)$

Keep $x_i$ with probability $\dfrac{p(x)}{K g(x)}$ where $K$

is a constant s.t. $p(x) \le K g(x) \; \forall x$

# Markov Chain Monte Carlo
## random sampling

Sequence of random variables where next step in
the sequence depends only on the previous

Efficient, multi-dim. sampling that preferentially

samples the "typical set."

Target distribution $\pi(\theta) = p(d|\theta, H) p(\theta|H)$

(switching notation a bit, $\pi$ is not a pdf because
it's not normalized)

proposal distribution $q(\theta^{(n+1)}|\theta^{(n)})$

acceptance prob $\alpha(\theta^{(n+1)}|\theta^{(n)})$

Transition prob. $P(\theta^{(n+1)}|\theta^{(n)})$ given by $q$ and $\alpha$

We want our chain to respect detailed balance
$$\pi(\theta^{(n+1)}) P(\theta^{(n)}|\theta^{(n+1)}) = \pi(\theta^{(n)}) P(\theta^{(n+1)}|\theta^{(n)})$$

i.e. probability of being in state $(n+1)$ and transitioning

to $(n)$ is the same as the reverse.

This yields a stationary distribution, i.e. if

$\theta^{(n)}$ drawn from $\pi$, $\theta^{(n+1)}$ also drawn from $\pi$

and we will "eventually" converge to the target

distribution

## Metropolis Hastings algorithm

$$\alpha(\theta^{(n+1)}|\theta^{(n)}) = \min\left[1, \frac{\pi(\theta^{(n+1)}) \, q(\theta^{(n)}|\theta^{(n+1)})}{\pi(\theta^{(n)}) \, q(\theta^{(n+1)}|\theta^{(n)})}\right]$$

In this case, transition prob is

$$p(\theta^{(n+1)}|\theta^{(n)}) = q(\theta^{(n+1)}|\theta^{(n)})\alpha(\theta^{(n+1)}|\theta^{(n)})$$
$$+ (1-b)\delta(\theta^{(n)} - \theta^{(n+1)})$$
$$b = \int d\theta^{(n+1)} \, q(\theta^{(n+1)}|\theta^{(n)})\alpha(\theta^{(n+1)}|\theta^{(n)})$$

### This choice satisfies detailed balance.

Proof: If $\theta^{(n+1)} = \theta^{(n)}$, trivial

If $\theta^{(n+1)} \neq \theta^{(n)}$,

$$\alpha(\theta^{n+1}|\theta^{n}) = \min\{1, r\}$$

If $r > 1$,

LHS: $\pi(\theta^n) \, q(\theta^{n+1}|\theta^n) \cdot 1$

RHS: $\pi(\theta^{n+1}) \, q(\theta^n|\theta^{n+1}) \cdot \dfrac{1}{r}$

$= \quad '' \qquad '' \qquad \cdot \dfrac{\pi(\theta^n) \, q(\theta^{n+1}|\theta^n)}{\pi(\theta^{n+1}) \, q(\theta^n|\theta^{n+1})}$ ✓

If $r < 1$,

LHS: $\pi(\theta^n) \, q(\theta^{n+1}|\theta^n) \cdot r$     similar cancellation

RHS: $\pi(\theta^{n+1}) \, q(\theta^n|\theta^{n+1}) \cdot 1$
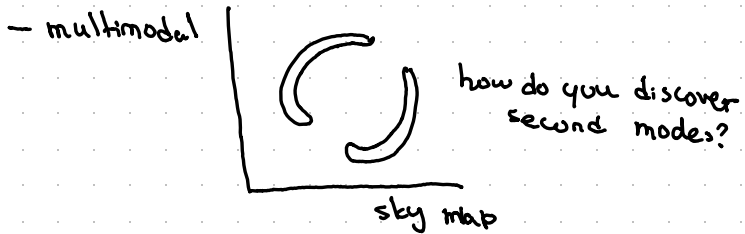
Note that often proposal $q$ is taken to be a Gaussian, so that

$$q(\theta^{n+1}|\theta^n) = q(\theta^n|\theta^{n+1})$$

and $r = \dfrac{\pi(\theta^{n+1})}{\pi(\theta^n)}$

# Challenges:

- correlations



good jump in one direction isn't necessarily
a good jump in another direction

- multimodal



how do you discover
second modes?

sky map
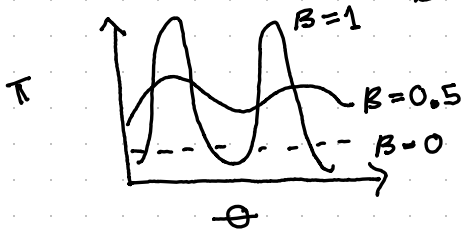
One technique is Parallel Tempering MCMC

Run multiple chains with different target
distributions

$$\pi_\beta(\theta) = \left[ p(d|\theta, H) \right]^\beta p(\theta|H)$$

$$\beta = \frac{1}{T} \leftarrow \text{temperature}$$



$\beta = 1$ (cold chain) is our posterior

$\beta \to 0$ (hot chain) is our prior

hot chain can hop around much more easily.

Progress the chains together
every N iterations, propose a swap between
neighboring chains i and j

accept the swap with probability

$$A_{ij} = \min\left[1, \frac{\pi_{B_i}(\theta_j)}{\pi_{B_j}(\theta_i)}\right]$$

$$= \min\left[1, \frac{p(d|\theta_j)}{p(d|\theta_i)}\right]^{B_i - B_j}$$

helps explore parameter space (find new modes)

also used to calculate evidences
("thermodynamic integration")

Define an evidence for each temperature chain

$$Z_B \equiv \int d\theta\, \pi_B(\theta) \quad \text{so that} \quad P_B(\theta) = \frac{\pi_B(\theta)}{Z_B}$$

[The evidence we <u>want</u> is $B=1$]

consider $\frac{\partial}{\partial B} \ln(Z_B) = \frac{1}{Z_B} \frac{\partial}{\partial B} Z_B$

$$= \frac{1}{Z_B} \int d\theta \frac{\partial}{\partial B} \pi_B(\theta)$$

$$= \int d\theta \frac{1}{Z_B} \frac{\partial \pi_B(\theta)}{\partial B} \cdot \frac{\pi_B(\theta)}{\pi_B(\theta)}$$

$$= \int d\theta\, P_B(\theta) \frac{\partial \ln \pi_B(\theta)}{\partial B}$$

$$\frac{\partial \ln \pi_B(\theta)}{\partial B} = \frac{\partial}{\partial B}\left[\ln\left[(p(d|\theta, H))^B\right] + \ln p(\theta|H)\right]$$

$$= \frac{\partial}{\partial B}\left[B \ln p(d|\theta, H)\right]$$

$$= \ln p(d|\theta, H) \quad \text{which is just the log-likelihood}$$

so $\frac{\partial}{\partial B} \ln(Z_B) = \int d\theta\, P_B(\theta) \ln p(d|\theta, H)$

$$= \mathbb{E}_B\left[\ln p(d|\theta, H)\right]$$

Easy to estimate for each chain $B$ by
taking average over points in the chain
Once we have $\mathbb{E}_B[\ln p(d|\theta, H)]$ computed for
each chain, we can numerically integrate

$$\int_0^1 \frac{\partial}{\partial B} \ln(Z_B)\, dB = \ln Z_1 - \ln Z_0$$

$$= Z_1 \qquad \text{"1 if prior is normalized}