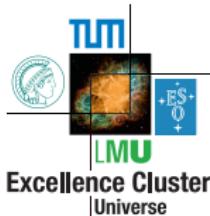


POPULATION MONTE CARLO

Frederik.Beaujean@lmu.de

Helmholtz school on
*Monte Carlo methods in advanced
statistics applications and data analysis*
Nov 22, 2013



PLAN

- ① MOTIVATION
- ② INTRODUCTION TO PMC
- ③ INITIALIZATION
- ④ COMPARISON AND OUTLOOK

SUMMARY

- population Monte Carlo (PMC) \equiv adaptive importance sampling
- complicated target $P \Rightarrow$ need flexible proposal function q

$$q(\vec{x}) = \sum_{j=1}^K \alpha_j q_j(\vec{x} \mid \vec{\mu}_j, \Sigma_j)$$

- adapt q iteratively to match P
- big issue: initialization

PLAN

① MOTIVATION

② INTRODUCTION TO PMC

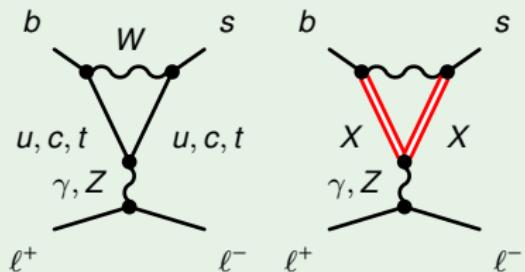
③ INITIALIZATION

④ COMPARISON AND OUTLOOK

NEW PHYSICS IN RARE B-MESON DECAYS

FLAVOR CHANGING NEUTRAL CURRENTS IN SM

- quark transition $b \rightarrow s$ at loop level
- Rare decay: $\mathcal{B}(B \rightarrow K\mu^+\mu^-) \sim 10^{-7}$
- indirect search for heavy particles up to $\mathcal{O}(\text{TeV})$

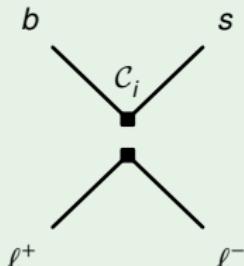


EFFECTIVE THEORY

- high energy: effective coupling (Wilson coefficient) C_i
- low energy: effective operator \mathcal{O}_i

$$\mathcal{L}_{\text{eff}} \sim \sum_i (\mathcal{C}_i^{\text{SM}} + \Delta C_i) \mathcal{O}_i$$

⇒ Extract \vec{C} in global fit



GLOBAL FIT

THE GOAL

- ① Compare SM with new physics:
Bayes factor
- ② Extract marginal distribution of
 \vec{C} , $\dim \vec{C} = 2 \dots 6$

BAYES' THEOREM

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

$$P(\vec{C}, \vec{\nu}|D) = \frac{P(D|\vec{C}, \vec{\nu})P(\vec{C})P(\vec{\nu})}{\int d\vec{C}d\vec{\nu} P(D|\vec{C}, \vec{\nu})P(\vec{C})P(\vec{\nu})}$$

$$\text{Bayes factor} = \frac{\text{evidence|SM}}{\text{evidence|NP}}$$

INPUTS

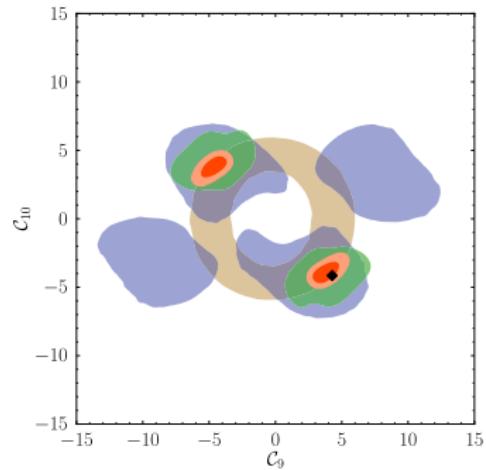
- theory uncertainty: quark masses, nonperturbative QCD ... $\Rightarrow P(\vec{\nu})$, $\dim \vec{\nu} \approx 30$
- over 50 results from Belle, LHCb ... experiments $\Rightarrow P(D|\vec{C}, \vec{\nu})$

COMBINATIONS OF INPUT

- ① all data:  1σ  2σ
- ② Standard Model: 
- ③ data subsets ...

CHALLENGES

- Multiple isolated maxima: observable $X \sim \mathcal{C}_i \mathcal{C}_j = (-\mathcal{C}_i) (-\mathcal{C}_j)$
- Theory prediction slow: ~ 1 s for one sample, but need $\mathcal{O}(10^6)$ samples \Rightarrow parallel evaluation
- 30 dimensions

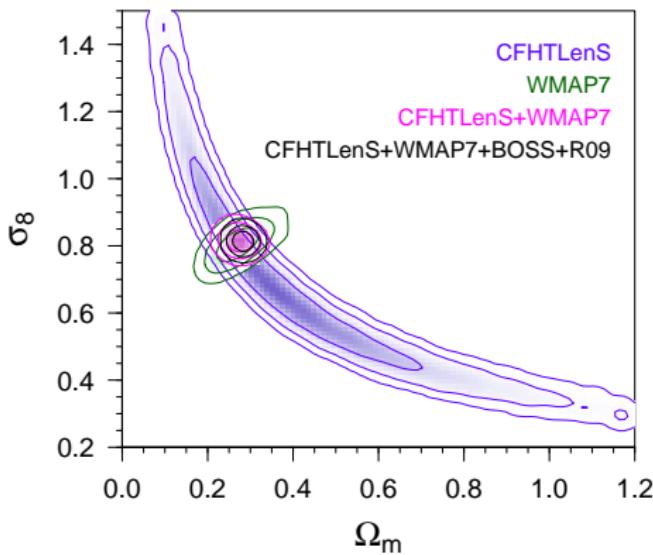


BAYES FACTOR

$$\frac{\text{evidence|SM}}{\text{evidence|NP}} = 200 \dots 1500$$

COSMOLOGY

- Kilbinger, Cappé, et al. develop **cosmoPMC** to constrain cosmological parameters in global fits
- **cosmoPMC** contains `pmclib`, C implementation of PMC with MPI parallelization



Kilbinger et al. arXiv:1212.3338

BLACK BOX — HOLY GRAIL OF MC SAMPLING

TARGET

- just code that returns the value at \vec{x}
- any number of modes, curvature or degeneracy
- any dimension

SAMPLER

- use any number of cores effectively
- no manual tuning required
- output always reliable

PLAN

① MOTIVATION

② INTRODUCTION TO PMC

③ INITIALIZATION

④ COMPARISON AND OUTLOOK

PMC BASICS I

Regular importance sampling: $\vec{x} \in V \subseteq \mathbb{R}^d$

$$I = \int d\vec{x} P(\vec{x}) = \int d\vec{x} \frac{P(\vec{x})}{q(\vec{x})} q(\vec{x}) \quad (1)$$

$$\Rightarrow \hat{I} = \frac{1}{N} \sum_{i=1}^N \frac{P(\vec{x}^i)}{q(\vec{x}^i)} = \frac{1}{N} \sum_{i=1}^N w^i, \vec{x} \sim q \quad (2)$$

Terminology:

q : proposal density (3)

P : target density (often: posterior) (4)

$w \equiv \frac{P}{q}$ importance weight (5)

PMC BASICS II

Self-normalized weights

$$\bar{w}^k \equiv \frac{w^k}{\sum_i w^i} \quad (6)$$

Good for arbitrary expectation

$$E_P[f] \equiv \int d\vec{x} P(\vec{x}) f(\vec{x}) \approx \sum_{i=1}^N \bar{w}^i f(\vec{x}^i) \quad (7)$$

Law of large numbers: converges as

$$\sim 1/\sqrt{N} \quad (8)$$

MARGINAL DISTRIBUTIONS

Main application: marginal distributions, example in 1D

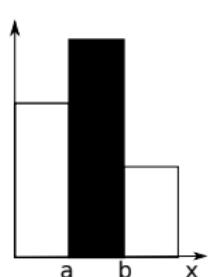
$$P(a \leq x_1 \leq b) = \int_a^b dx_1 \int dx_j P(\vec{x}) \quad (9)$$

$$P(a \leq x_1 \leq b) \approx \sum_{i=1}^N \bar{w}^i f(\vec{x}^i) \quad (10)$$

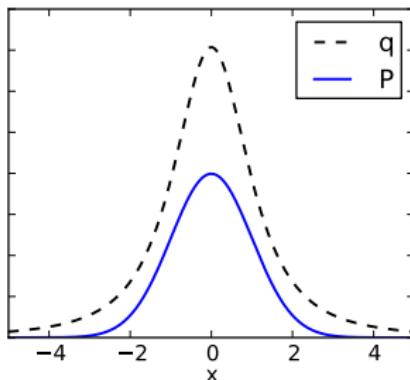
with the indicator function

$$f(\vec{x}) = \begin{cases} 1, & x_1 \in [a, b] \\ 0, & \text{else} \end{cases} \quad (11)$$

Note: For Markov chains, $\vec{x} \sim P \Rightarrow w \equiv 1$



PROPOSAL DENSITY



Requirements:

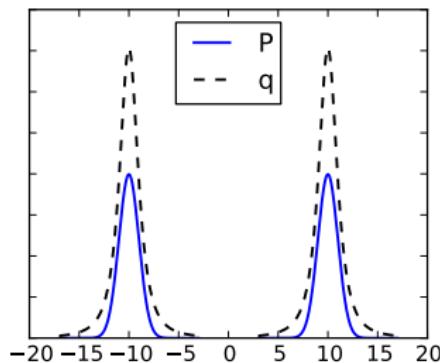
- $\text{supp } P \subseteq \text{supp } q$
- normalized
 \Rightarrow get normalization of P
- falls off slower than P for large \bar{x}
 \Rightarrow estimators have finite variance

Desirable:

- easy to draw from
- for most applications:

$$q \propto P \text{ optimal}$$

MIXTURE DENSITY



A very flexible class of functions:

$$q(\vec{x}) = \sum_{j=1}^K \alpha_j q_j(\vec{x} | \vec{\mu}_j, \Sigma_j) \quad (12)$$

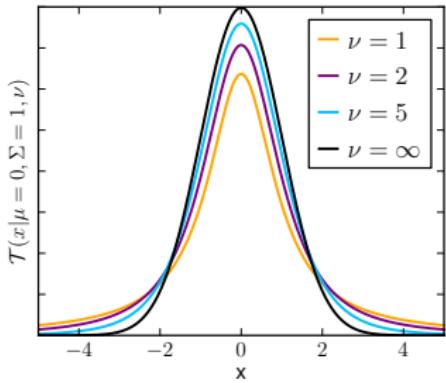
q_j = Gauss, \mathcal{N} or Student's t distribution, \mathcal{T}

α_j = component weight, $\sum_j \alpha_j = 1$

STUDENT'S T

$$\mathcal{T}(\vec{x}|\vec{\mu}, \Sigma, \nu) = \frac{\Gamma((\nu + d)/2)}{\Gamma(\nu/2)(\pi\nu)^{d/2}} |\Sigma|^{-1/2} \left(1 + \frac{1}{\nu} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right)^{-(\nu+p)/2} \quad (13)$$

- mean $\vec{\mu}$, covariance $\Sigma \frac{\nu}{\nu - 2}$
- $\nu \equiv$ degree of freedom
- $\nu = 1 \Rightarrow$ Cauchy, Lorentz, Breit-Wigner
- $\nu = \infty \Rightarrow$ Gauss, \mathcal{N}



wikipedia

HOW TO ADAPT q ?

1. need to measure discrepancy between P, q
2. then minimize discrepancy

HOW TO ADAPT q ?

1. need to measure discrepancy between P, q
 2. then minimize discrepancy
-
1. Kullback-Leibler divergence (from P to q)

$$\text{KL}(P\|q) = \int d\vec{x} P(\vec{x}) \log \frac{P(\vec{x})}{q(\vec{x})} \quad (14)$$

- $\text{KL}(P\|q) \neq \text{KL}(q\|P) \Rightarrow \text{KL}$ not a metric
- $\text{KL} \geq 0$ and $\text{KL} = 0 \Leftrightarrow P \equiv q$
- KL hard to compute, even harder than $\int d\vec{x} P(\vec{x})$

HOW TO ADAPT q ?

1. need to measure discrepancy between P, q
2. then minimize discrepancy
2. Simplification: restrict to mixture density $q = \sum_{j=1}^K \alpha_j q_j(\cdot | \mu_j, \Sigma_j)$
 - functional optimization \rightarrow ordinary parameter optimization
 - Global fit [arXiv:1205.1838](https://arxiv.org/abs/1205.1838)

$$d = 30, K = 100 \Rightarrow \dim \alpha + (\dim \vec{\mu} + \dim \Sigma) \times \dim \alpha \quad (14)$$

$$100 + (30 + 30 \times 31/2) \times 100 \quad (15)$$

$$= 49600 \text{ parameters to fix} \quad (16)$$

- too much for gradient-based methods like Minuit!
- can't even compute $\text{KL}(P||q)$!

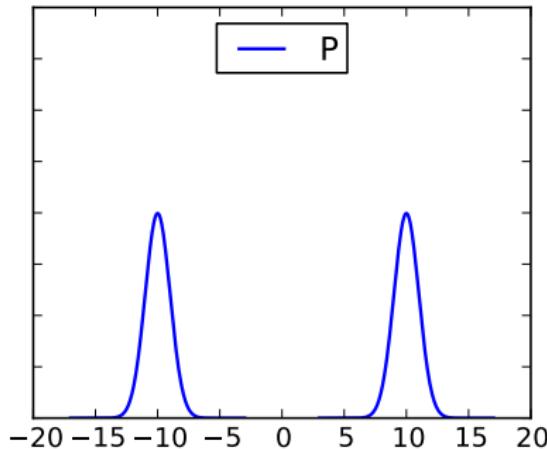
A CLEVER SOLUTION

- use analytical structure of the problem
- and expectation maximization (EM) Dempster, Laird, Rubin (1977)
- easy reading: Piater (2002) Borman (2004)

A CLEVER SOLUTION

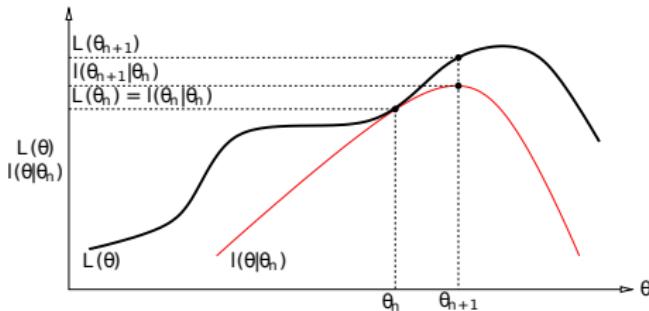
- use analytical structure of the problem
- and expectation maximization (EM) [Dempster, Laird, Rubin \(1977\)](#)
- easy reading: [Piater \(2002\)](#) [Borman \(2004\)](#)

EM — BASIC IDEA



- Observe $D = \{x^i : i = 1 \dots N\}$
- Assume $x^i \sim P(\cdot) = \sum_{j=1}^K \alpha_j \mathcal{N}(\cdot | \mu_j, \sigma_j)$
- Maximum-likelihood fit of mixture to data
- If we knew generating component $z_i \in \{1, \dots, K\} \Rightarrow K$ indep. problems
- Solution easy: $\mu_j = 1/N_j \sum_i x^i$, $\sigma_j^2 = 1/N_j^2 \sum_i (x^i - \mu_j)^2$

EXPECTATION MAXIMIZATION IN A NUTSHELL



Borman (2004)

- iterative procedure to maximize L : E step \leftrightarrow M step
- E step: infer missing data given current θ_n and observed data
- M step: maximize simpler function: $\theta_{n+1} = \arg \max I(\theta|\theta_n)$
- converge to *local* optimum in $\mathcal{O}(10)$ steps

- ① Initial guess q^0 at $t = 0$
- ② Draw N samples $\vec{x}^i \sim q^t$, compute $w^i = \frac{P(\vec{x}^i)}{q(\vec{x}^i)}$
- ③ Update q^t based on $\{(\vec{x}^i, w^i) : i = 1, \dots, N\}$
- ④ Repeat steps 2 and 3 until convergence
- ⑤ If converged, regular IS with q^t

NOTE

- computing weights usually slowest part
- samples in step 2 indep. \Rightarrow **massive parallelization**
- samples form *population* \Rightarrow population Monte Carlo

CONVERGENCE — OR WHEN TO STOP

Optimization goal

$$\text{KL}(P\|q^t) = 0 \Rightarrow \exp[\text{KL}(P\|q^t)] = 1 \quad (17)$$

estimated by *normalized perplexity*

$$\mathcal{P}^t \equiv \frac{\exp(H^t)}{N} \in [0, 1] \quad (18)$$

Shannon entropy of self-normalized weights [Shannon \(1948\)](#)

$$H(\{\bar{w}_i^t\}) = - \sum_{i=1}^N \bar{w}_i^t \log \bar{w}_i^t \quad (19)$$

Compare with quantum mechanics

$$S = -k_B \sum_i p_i \log p_i, \text{ where } p_i \equiv P(\text{state } i) \quad (20)$$

Maximum entropy for

$$q^t = P \Rightarrow \bar{w} = 1/N \Rightarrow H^t = \log N \quad (21)$$

SAMPLE QUALITY

What fraction of samples actually contributes? *effective sample size*

$$\text{ESS}^t \equiv \left(\sum_{i=1}^N \{\bar{w}_i^t\}^2 \right)^{-1} / N \in [1/N, 1] \quad (22)$$

Examples:

$$\bar{w} = 1/N \Leftrightarrow \text{ESS} = 1 \quad (23)$$

$$\bar{w} = \underbrace{\{0, \dots, 0\}}_{N_0} \underbrace{\{c, \dots, c\}}_{N-N_0} \Rightarrow \text{ESS} = \frac{N - N_0}{N} \quad (24)$$

Extreme case: single outlier dominates all samples

$$N_0 = N - 1 \Rightarrow \text{ESS} = 1/N \quad (25)$$

CONVERGENCE CRITERION

$$\mathcal{P}^t = \mathcal{P}^t \left(\overline{w}_i^t \log \overline{w}_i^t \right)$$

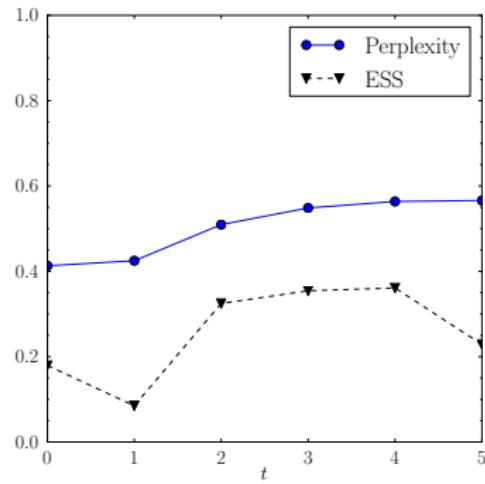
vs

$$\text{ESS}^t = \text{ESS}^t \left(\{\overline{w}_i^t\}^2 \right)$$

- ESS better indicator of outliers
- \mathcal{P} more robust indicator
- stop adapting q^t if in last two steps

$$\frac{\sqrt{V[\mathcal{P}]} }{E[\mathcal{P}]} \leq \epsilon \quad (26)$$

with $\epsilon \sim 10\%$



A tough example in 22D

How to start?

PLAN

① MOTIVATION

② INTRODUCTION TO PMC

③ INITIALIZATION

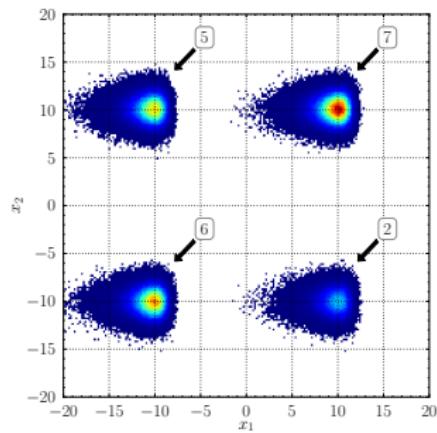
④ COMPARISON AND OUTLOOK

MONTE CARLO SAMPLING

Goal: draw samples from target P

1 Markov chain Monte Carlo (MCMC)

	Pro	Con
MCMC	local exploration, learns on the fly	trapped in local maxima



FOUR IDENTICAL MAXIMA

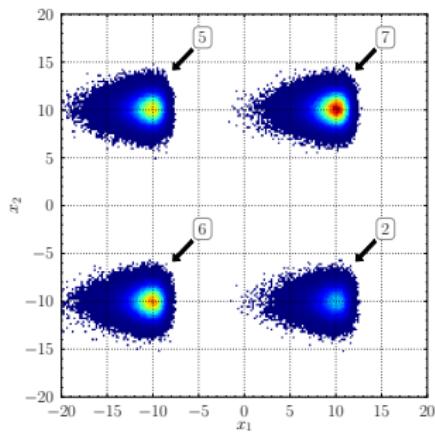
- start 20 local-random-walk chains
- each trapped in one mode

MONTE CARLO SAMPLING

Goal: draw samples from target P

- ① Markov chain Monte Carlo (MCMC)
- ② adaptive importance sampling = population Monte Carlo (PMC)

	Pro	Con
MCMC	local exploration, learns on the fly	trapped in local maxima
PMC	massive parallelization, yields normalization, multiple modes OK	very sensitive to initialization



FOUR IDENTICAL MAXIMA

- start 20 local-random-walk chains
- each trapped in one mode

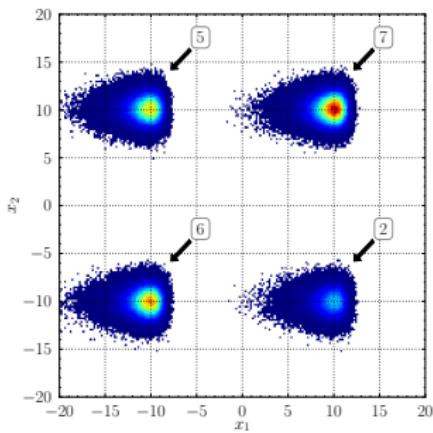
MONTE CARLO SAMPLING

Goal: draw samples from target P

- ① Markov chain Monte Carlo (MCMC)
- ② adaptive importance sampling = population Monte Carlo (PMC)

	Pro	Con
MCMC	local exploration, learns on the fly	trapped in local maxima
PMC	massive parallelization, yields normalization, multiple modes OK	very sensitive to initialization

→ Combine MCMC and PMC

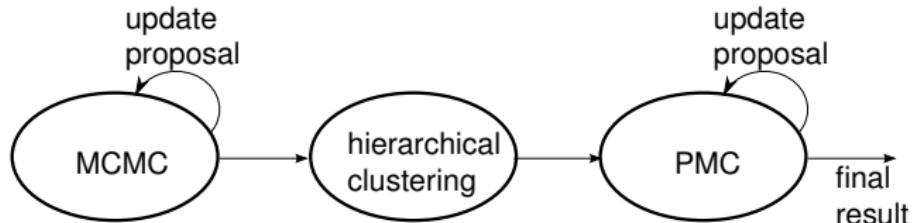


FOUR IDENTICAL MAXIMA

- start 20 local-random-walk chains
- each trapped in one mode

INITIALIZING PMC

BEAUJEAN & CALDWELL (2013)

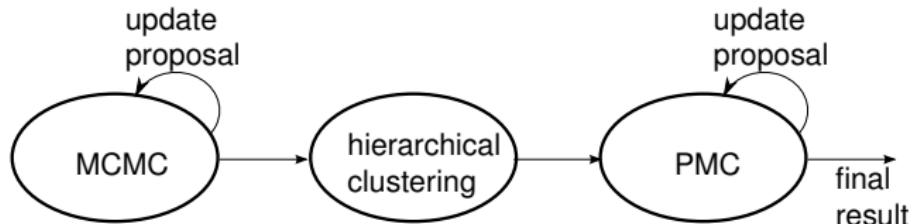


MARKOV CHAIN IMPLEMENTATION HAARIO ET AL. (2001)

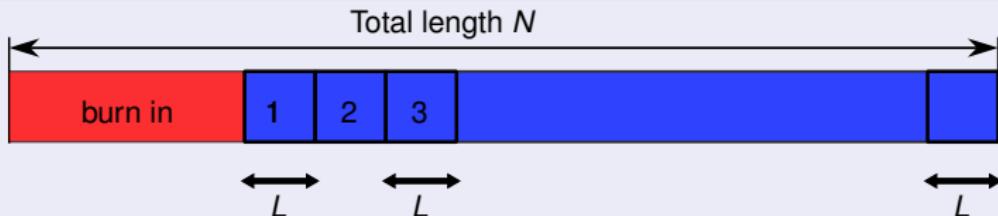
- starting point drawn uniformly across parameter space
- Gaussian proposal whose covariance is
 - adapted to target, $\Sigma \rightarrow \Sigma_P$
 - scaled such that proposed points are accepted with (15 – 35)%.
- don't have to stop adapting, quit after N steps
- chains stuck \Rightarrow multiple chains

INITIALIZING PMC

BEAUJEAN & CALDWELL (2013)



SPLIT UP THE LOCAL-RANDOM-WALK MARKOV CHAIN

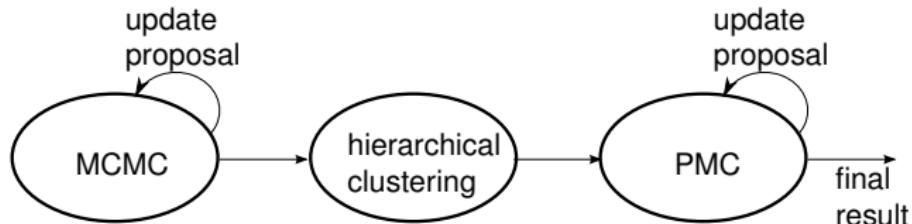


patch of length $L \stackrel{!}{=} \text{one component}$

$$\text{sample mean } \bar{\mu} = \frac{1}{N} \sum_{i=1}^L \bar{x}^i, \quad \text{sample cov. } \Sigma = \frac{1}{N-1} \sum_{i=1}^L (\bar{x}^i - \bar{\mu})(\bar{x}^i - \bar{\mu})^T$$

INITIALIZING PMC

BEAUJEAN & CALDWELL (2013)



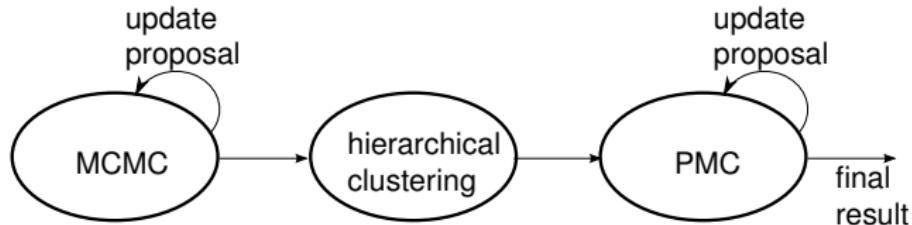
HIERARCHICAL CLUSTERING I

Goal: condense information

- Have Gaussian mixture with M components $f(\vec{x}) = \sum_{l=1}^M \beta_l f_l(\vec{x} | \vec{\mu}_l, \Sigma_l)$,
- Ex: 20 chains, 10000 samples, 20% burn-in, $L = 100 \Rightarrow M = 1600$ (**huge!**)
- Want mixture with $K \ll M$ components $q(\vec{x}) = \sum_{j=1}^m \alpha_j q_j(\vec{x} | \vec{\mu}_j, \Sigma_j)$, say $K = 10$
- Find q “closest” to f . Sounds familiar?

INITIALIZING PMC

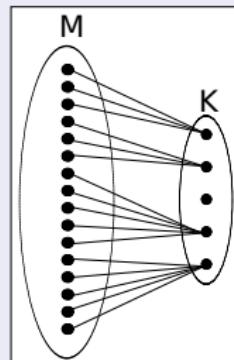
BEAUJEAN & CALDWELL (2013)



HIERARCHICAL CLUSTERING II

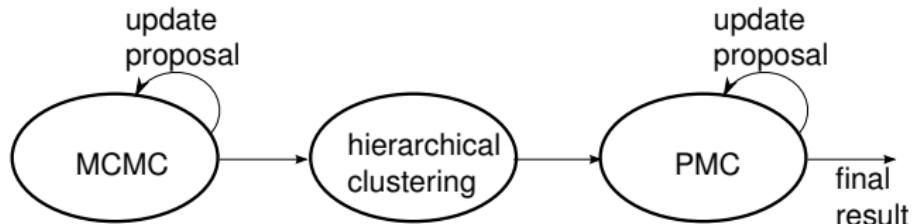
GOLDBERGER & ROWEIS (2004)

- distance measure $d(f, q) = \sum_{l=1}^M \alpha_l \min_j \text{KL}(f_l \| q_j)$,
- $\text{KL}(f_l \| q_j)$ simple for Gaussians \Rightarrow fast
- map many f_l to closest q_j , adjust q_j to minimize $d(f, q)$
- use expectation-maximization, need initial guess again



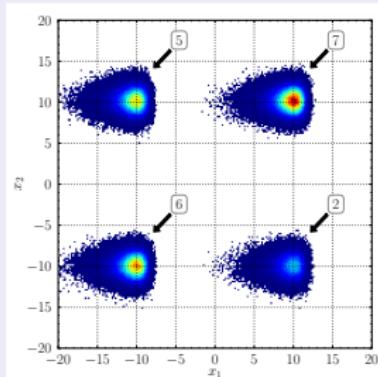
INITIALIZING PMC

BEAUJEAN & CALDWELL (2013)

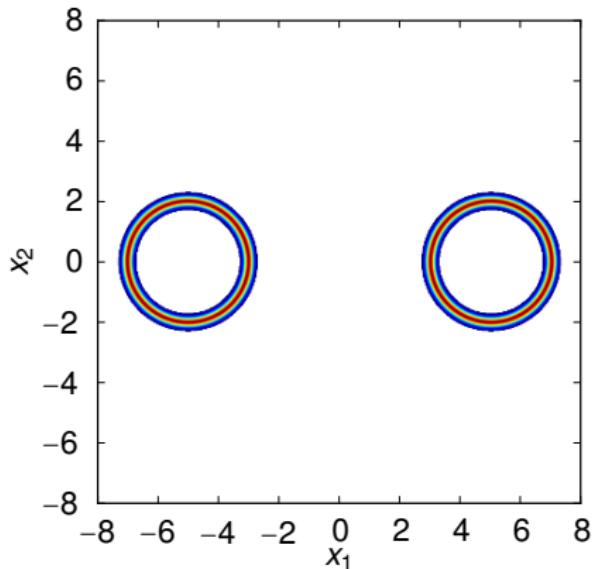


HIERARCHICAL CLUSTERING III: INITIAL GUESS

- pick K (too small \Rightarrow poor proposal, too large \Rightarrow wasting time)
- group chains: need similar mean and variance in *each* dimension
- same #components from each group
- sample mean, covariance from *long* chain patches



EXAMPLE: GAUSSIAN SHELLS IN 2D



$$P(\vec{x}) = \frac{1}{2} \text{circ}(\vec{x}|\vec{c}_1, r, w) + \frac{1}{2} \text{circ}(\vec{x}|\vec{c}_2, r, w)$$

$$\text{circ}(\vec{x}|\vec{c}, r, w) = \frac{1}{\sqrt{2\pi w^2}} \exp\left[-\frac{(|\vec{x} - \vec{c}| - r)^2}{2w^2}\right]$$

$$\vec{c} = (\pm 5, 0), r = 2, w = 0.1$$

CHALLENGES

- multimodal
- degeneracy

GETTING YOUR HANDS DIRTY

GET YOUR COPY OF THE EXERCISES

```
ssh -X USER@naf-school01.desy.de
% source /afs/desy.de/group/school/mc-school/setup_pmc.sh
% cp /afs/desy.de/user/s/school09/exercises.tar.bz ~
% tar xaf exercises.tar.bz ; ls exercises/
clustering mcmc pmc
```

- each folder has `exercise.pdf` with instructions
- use `evince` to open pdf files
- start with `mcmc`

PLAN

- ① MOTIVATION
- ② INTRODUCTION TO PMC
- ③ INITIALIZATION
- ④ COMPARISON AND OUTLOOK

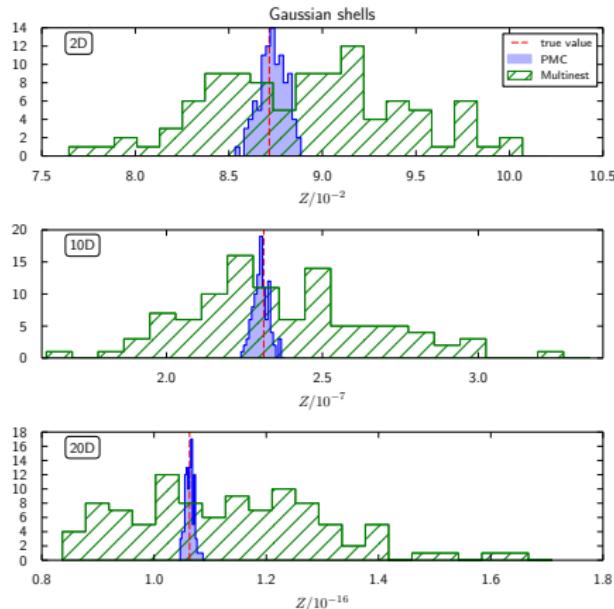
STATISTICAL UNCERTAINTY

$$\widehat{Z} = \frac{1}{N} \sum_{i=1}^N w^i$$

$$\widehat{\Delta Z} = \sqrt{\frac{1}{N(N-1)} \sum_{i=1}^N (w^i - \widehat{Z})^2}$$

PMC vs MULTINEST v2.18

- “same number” of samples contribute
- $d = 2 : \widehat{\Delta Z} \lesssim 1\%$ (PMC) vs 5% (nest)
- variance of \widehat{Z} in repetitions matches $\widehat{\Delta Z}$ very well



100 runs on Gaussian shells

IMPLEMENTATION

CURRENT STATUS

- Exercises use flavor tool `eos` ([docs](#)) and `pmclib` (backend)
- `matplotlib` (plots) and `HDF5` (I/O)

FUTURE DEVELOPMENT

- Publish MCMC + clustering + PMC in python ⇒
<https://github.com/fredRos/pypmc.git>

IMPROVEMENTS

GOALS

- User friendliness: less (no?) tuning parameters
- more robust results

INITIALIZATION

- Hierarchical clustering → variational Bayes
- sampling from Beta distr. ⇒ **massive parallelization**

PMC

- update student's t ν
- combine samples from multiple steps [Cornuet et al.](#) ⇒ soften outlier

$$w_t^i = P(\vec{x}^i) \Bigg/ \sum_t q^t(\vec{x}^i)$$

CONCLUSION

PMC

- is a powerful algorithm easily using 1000's of cores
- benefits from MCMC for initialization
- is not the holy grail