

Advanced Topics in Machine Learning 7 Inference - Part 2

Prof. Dr. Steffen Staab Dr. Rafika Boutalbi Zihao Wang https://www.ipvs.uni-stuttgart.de/departments/ac/









Learning Objectives

Inference Methods:

- MAP: Max-Sum
- Sampling
 - Forward sampling
 - MCMC
 - Gibbs Sampling
 - Metropolis Hasting

Disclaimer

Figures and examples not marked otherwise are taken from the book by Koller & Friedman

1 Max-Sum Exact Inference for Cluster Trees

Turn Product into Sum

$$P_{\Phi}(\mathbf{x}) \propto \prod_{k} \phi_{k} (D_{k} = \mathbf{x}_{k})$$

$$\operatorname{argmax} \ \prod_{k} \phi_{k} \big(D_{k} \big)$$

$$\operatorname{argmax} \sum_{k} \theta_{k}(D_{k}) = \operatorname{argmax} \theta(\boldsymbol{\mathcal{X}}),$$

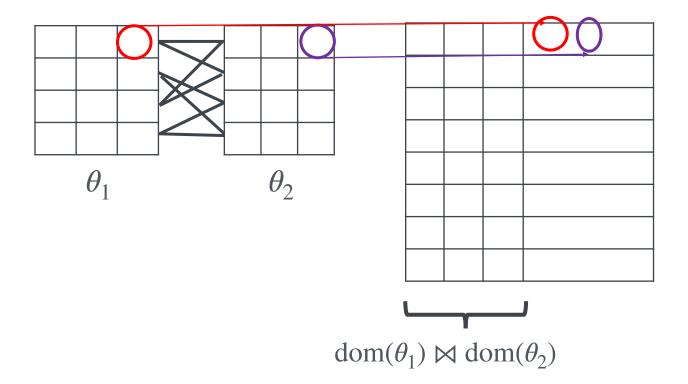
where
$$\theta_k(D_k) = \log \phi_k(D_k)$$

Remember, e.g., Naive Bayes in ML lecture!

Factor Summation

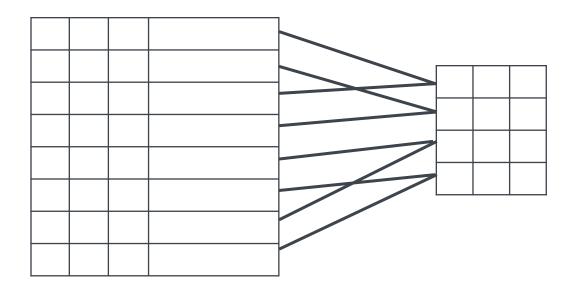
Example

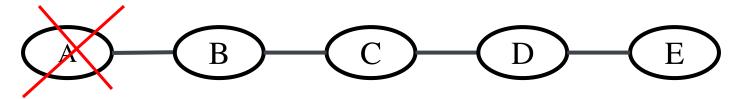
$$\theta_1 + \theta_2$$



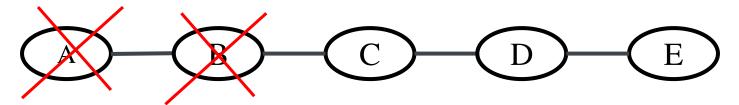
Max Marginalization

$$\begin{array}{l} \theta_1 + \theta_2 \\ b \in B \end{array} \qquad \text{argmax}(\theta_1 + \theta_2)$$





$$\max_{D}\max_{C}\max_{B}\max_{A}\left(\theta_{1}(A,B)+\theta_{2}(B,C)+\theta_{3}(C,D)+\theta_{4}(D,E)\right) = \\ =\max_{D}\max_{C}\max_{B}\left(\theta_{2}(B,C)+\theta_{3}(C,D)+\theta_{4}(D,E)+\\ \\ =\max_{D}\max_{C}\max_{B}\left(\theta_{2}(B,C)+\theta_{3}(C,D)+\theta_{4}(D,E)+\\ \\ =\max_{D}\max_{C}\max_{B}\left(\theta_{2}(B,C)+\theta_{3}(C,D)+\theta_{4}(D,E)+\\ \\ =\max_{D}\max_{C}\max_{B}\left(\theta_{2}(B,C)+\theta_{3}(C,D)+\\ \\ =\max_{D}\left(\theta_{2}(B,C)+\theta_{3}(C,D)+\\ \\ =\max_{D}\left(\theta_{2}(B,C)+\theta_{3}(C,D)+\\ \\ =\max_{D}\left(\theta_{2}(B,C)+\\ \\ =\max_{D}\left(\theta_{2}(B,$$



 $\max_{D} \max_{C} \max_{B} \max_{A} \left(\theta_1(A,B) + \theta_2(B,C) + \theta_3(C,D) + \theta_4(D,E) \right) = \max_{D} \max_{C} \left(\theta_2(B,C) + \theta_3(C,D) + \theta_4(D,E) + \max_{A} \left(\theta_2(B,C) + \theta_3(C,D) + \theta_4(D,E) + \max_{B} \left(\theta_2(B,C) + \theta_4(D,E) + \max_{B} \left(\theta_2(B,C) + \theta_4(D,E) + \max_{B} \left(\theta_2(B,C) + \theta_4(D,E) + \alpha_B \right) + \alpha_B \left(\theta_2(B,C) + \alpha_B \right) \right) \right) \right) \right)$



 $\max_{D}\max_{C}\max_{B}\left(\theta_{1}(A,B)+\theta_{2}(B,C)+\theta_{3}(C,D)+\theta_{4}(D,E)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\theta_{3}(C,D)+\theta_{4}(D,E)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\theta_{3}(D,E)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\theta_{2}(D,E)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\theta_{2}(B,C)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\theta_{2}(B,C)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\theta_{2}(B,C)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\theta_{2}(B,C)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\theta_{2}(B,C)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\theta_{2}(B,C)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\theta_{2}(B,C)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\lambda_{1}(B)\right)=\\ =$



 $\max_{D}\max_{C}\max_{B}\max_{A}\left(\theta_{1}(A,B)+\theta_{2}(B,C)+\theta_{3}(C,D)+\theta_{4}(D,E)\right)=\\ =\max_{D}\max_{C}\max_{B}\left(\theta_{2}(B,C)+\theta_{3}(C,D)+\theta_{4}(D,E)+\lambda_{1}(B)\right)=\\ =\max_{D}\max_{C}\left(\theta_{3}(C,D)+\theta_{4}(D,E)+\lambda_{2}(C)\right)=\\ =\max_{D}\max_{C}\left(\theta_{2}(B,C)+\theta_{3}(C,D)+\theta_{4}(D,E)+\lambda_{2}(C)\right)=\\ =\max_{D}\max_{C}\left(\theta_{1}(B,C)+\theta_{2}(C,D)+\theta_{4}(D,E)+\lambda_{2}(C)\right)=\\ =\max_{D}\max_{C}\left(\theta_{1}(B,C)+\theta_{2}(C,D)+\theta_{4}(D,E)+\lambda_{2}(C)\right)=\\ =\max_{D}\max_{C}\left(\theta_{1}(B,C)+\theta_{2}(C,D)+\theta_{4}(D,E)+\lambda_{2}(C)\right)=\\ =\max_{D}\max_{C}\left(\theta_{1}(B,C)+\theta_{2}(C,D)+\theta_{4}(D,E)+\lambda_{2}(C)\right)=\\ =\max_{D}\max_{C}\left(\theta_{1}(B,C)+\theta_{2}(C,D)+\theta_{4}(D,E)+\lambda_{2}(C)\right)=\\ =\max_{D}\left(\theta_{1}(B,C)+\theta_{2}(D,E)+\lambda_{2}(C)\right)=\\ =\max_{D}\left(\theta_{1}(B,C)+\theta_{2}(D,E)+\lambda_{2}(C)\right)=\\ =\max_{D}\left(\theta_{2}(B,C)+\theta_{3}(C,D)+\theta_{4}(D,E)+\lambda_{2}(C)\right)=\\ =\max_{D}\left(\theta_{1}(B,C)+\theta_{2}(D,E)+\lambda_{2}(C)\right)=\\ =\max_{D}\left(\theta_{1}(B,C)+\theta_{2}(D,E)+\lambda_{2}(D)+\theta_{2}(D,E)+\lambda_{2}(D)+\theta_{2}(D,E)+\lambda_{2}(D)+\theta_{2$

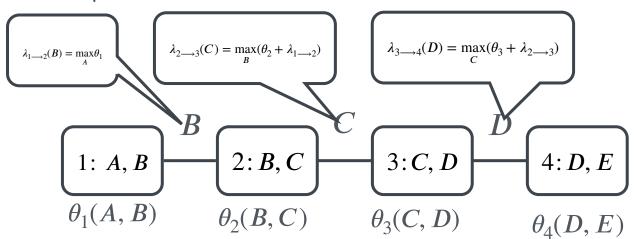
$$=\lambda_4(e)$$

Max-Sum in Cluster Trees

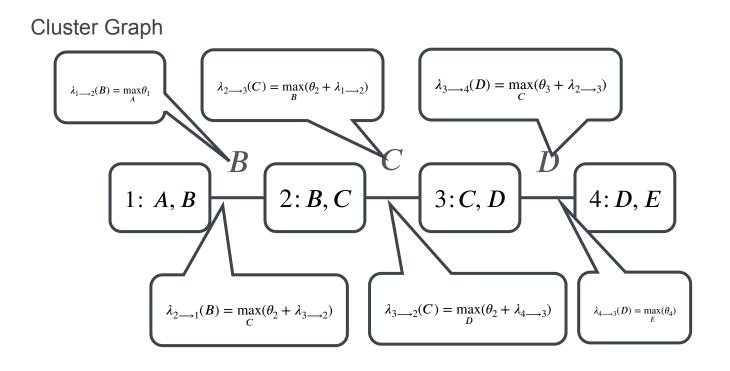
Factor Graph



Cluster Graph



Max-Sum in Cluster Trees



$$\theta_1(A, B)$$

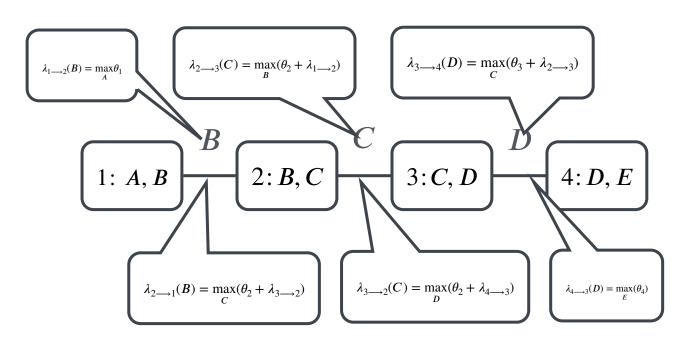
$$\theta_2(B,C)$$

$$\theta_2(B,C)$$
 $\theta_3(C,D)$

$$\theta_4(D, E)$$

Convergence of Message Passing

Once cluster C_i receives a final message from all neighbors except C_i then $\lambda_{i \longrightarrow i}$ is also final



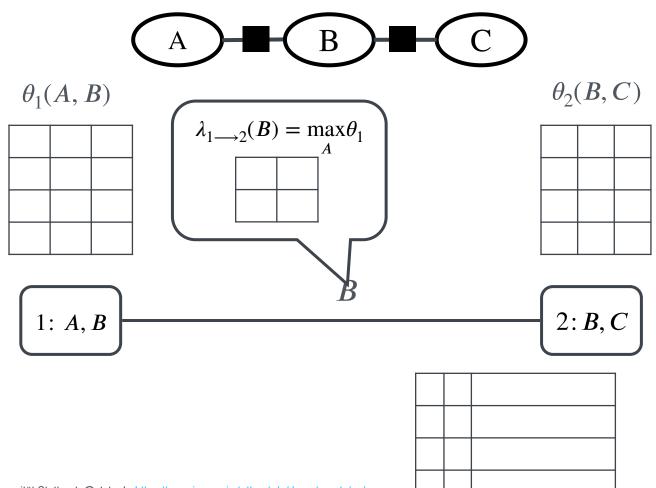
$$\theta_1(A, B)$$

$$\theta_2(B,C)$$

$$\theta_2(B,C)$$
 $\theta_3(C,D)$

$$\theta_4(D, E)$$

Example from before (upward pass)



Max-Sum Beliefs at Convergence

Beliefs at each Cluster are max-marginals

$$\beta_i(C_i) = \theta_i(C_i) + \sum_k \lambda_{k \longrightarrow i}$$

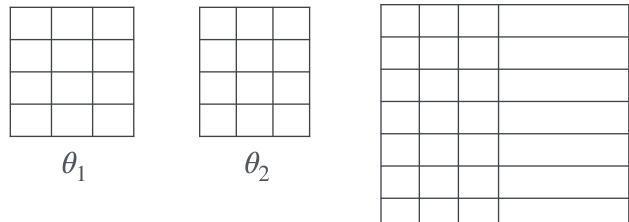
$$\beta_i(C_i) = \max_{W_i} \theta(C_i, W_i), \text{ where } W_i = \{X_1, ..., X_n\} - C_i$$

Calibration: clusters agree on shared variables

$$\max_{C_i - S_{i,j}} \beta_i(C_i) = \max_{C_j - S_{i,j}} \beta_j(C_j)$$

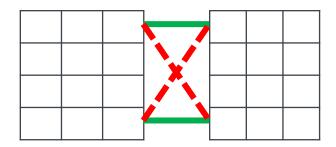
From Max to Argmax: Decoding a MAP Assignment

- Easy if MAP assignment is unique
 - Single maximizing assignment at each cluster
 - whose value is the θ value of the MAP assignment
 - Due to calibration, choices at all clusters must agree



Decoding Non-unique MAP assignment

XOR Problem

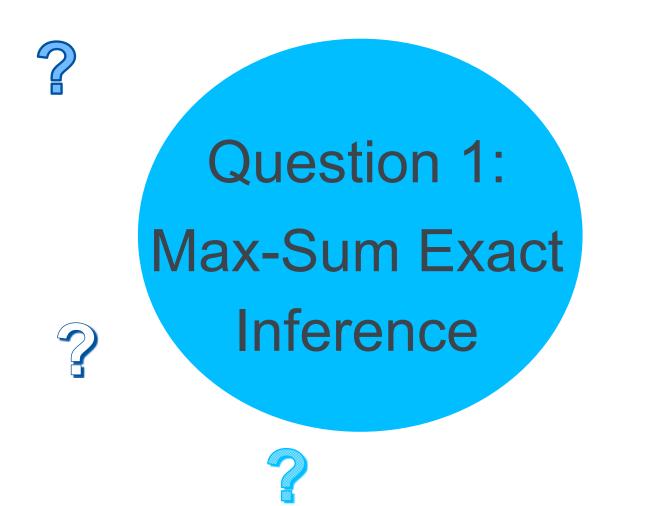


- If MAP assignment is not unique there are multiple choices at some clusters
- Arbitrary tie-breaking may not produce a MAP assignment
- Option 1: Slightly perturb parameters to make MAP unique
- Option 2: Use traceback that builds MAP assignment incrementally, one variable at a time

Summary

- The same cluster tree algorithm used for sum-product can be used for max-sum
- In cluster trees, convergence is achieved after a single updown pass
- Result is max-marginal at each cluster

There are tractable algorithms for MAP queries for many models if the models have special properties









2 Forward Sampling

Sampling-Based Estimation

- Data set $\mathcal{D} = \{x[1], ...x[M]\}$ sampled iid (independent, identically distributed) from $P(\mathcal{X})$
- Let $p := P(X = 1), X = \{0,1\}$
- Maximum likelihood estimation for p: $p_{\mathcal{D}} = \frac{1}{M} \sum_{m=1}^{M} x[m]$
- Expectation for a function f:

$$\mathbb{E}_{P(\mathcal{X})}[f] \approx \frac{1}{M} \sum_{m=1}^{M} f(x[m])$$

• If f is an indicator function, e.g. $f_0(x) = [x = 0]$ or $f_1(x) = [x = 1]$, we can use this to count occurrences

Sampling from a Discrete Distribution

$$\operatorname{Val}(\boldsymbol{\mathcal{X}}) = \{x^1, \dots, x^k\}, \quad P(x^i) = \theta^i$$

Python: Return random float uniformly from range [0.0, 1.0). random.random()

$$\begin{array}{ccc}
\theta^1 + \theta^2 \\
0 & \theta^1 \\
x^1 & x^2
\end{array}
\qquad \qquad \theta^1 + \theta^2 + \theta^3 \qquad \cdots 1$$

$$\theta^{1} + \theta^{2} + \theta^{3} + \dots + \theta^{k} = 1$$

Quality of Sampling-Based Estimation

Hoeffding Bound:

$$P_{\mathfrak{D}}(p_{\mathfrak{D}} \notin [p-\epsilon, p+\epsilon]) \le 2e^{-2M\epsilon^2}$$

• For additive bound ϵ on error with probability $> 1 - \delta$:

$$M \ge \frac{\ln\left(\frac{2}{\delta}\right)}{2\epsilon^2}$$

Chernoff Bound:

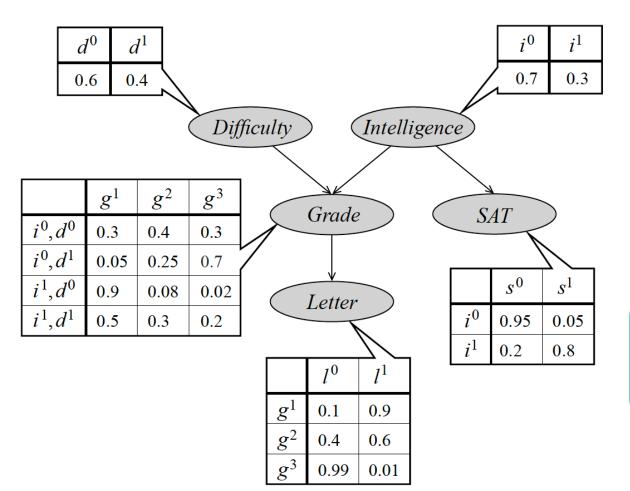
$$P_{\mathfrak{D}}(p_{\mathfrak{D}} \notin [p(1-\epsilon), p(1+\epsilon)]) \le 2e^{-Mp\epsilon^2/3}$$

• For multiplicative bound ϵ on error with probability $> 1 - \delta$:

$$M \ge 3 \frac{\ln\left(\frac{2}{\delta}\right)}{p\epsilon^2}$$

In other words:
sampling is
problematic when
p is small

Forward sampling from a Bayesian Network



 d^{0} , i^{1} , g^{1} , s^{0} , l^{1} d^{1} , i^{1} , g^{2} , s^{1} , l^{0}

Sample individual variables starting with parents before children

Querying by Forward Sampling

- Goal: Estimate P(Y = y)
 - 1. Generate M samples from the Bayesian network
 - 2. Compute fraction where Y = y
- . For additive bound ϵ on error with probability $>1-\delta$: $M\geq \frac{\ln\left(\frac{2}{\delta}\right)}{2\epsilon^2}$
- . For multiplicative bound ϵ on error with probability $>1-\delta$: $M\geq 3\frac{\ln\left(\frac{2}{\delta}\right)}{p\epsilon^2}$

Querying with Evidence by Rejection Sampling

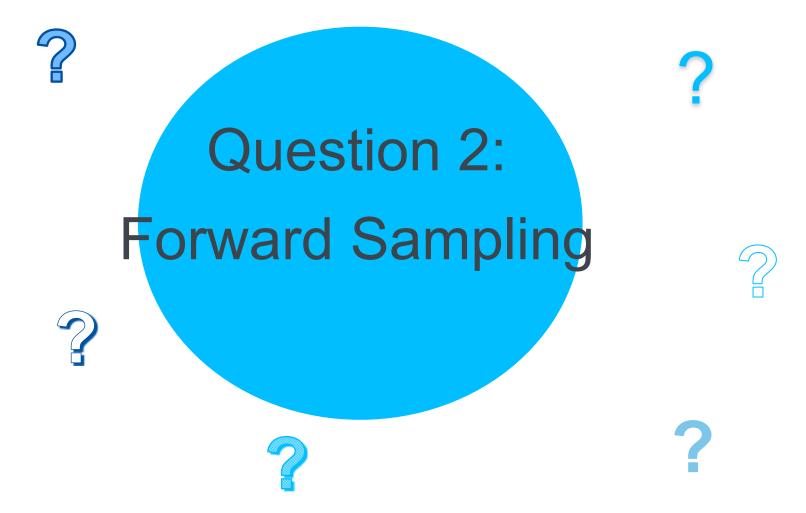
- Goal: Estimate P(Y = y | E = e)
 - 1. Generate samples from the Bayesian network
 - 2. Reject all samples where $E \neq e$
 - 3. Compute fraction of remaining samples where Y = y
- Expected fraction of samples kept: P(e)
- Curse of dimensionality:

Rejection sampling methods suffer from the *curse of dimensionality*, where the probability of rejection increases exponentially as a function of the number of dimensions |E|.

Summary

- Generating samples from a Bayesian Network is easy
- (ϵ, δ) -bounds exist, but usefulness is limited:
 - additive bounds: useless for low probability events
 - multiplicative bounds: # of samples grows as $\frac{1}{P(y)}$
- With evidence, # of required samples grows exponentially with # of observed variables

Forward sampling generally infeasible for Markov Networks.



3 Markov Chain

Markov Chain Monte Carlo (MCMC)

Markov chain Monte Carlo (MCMC) methods are a class of algorithms for *sampling from a probability distribution*.

They are **based on a Markov chain** that has the desired distribution of its **equilibrium** distribution.

The state of the chain after a number of steps is then used for sampling from the desired distribution.

MCMC algorithms are generally used for sampling from multi-dimensional distributions, especially when the *number of dimensions is high*.

Variations:

- Gibbs sampling
- Metropolis Hastings

Markov Chain Monte Carlo (MCMC)

Markov chain Monte Carlo (MCMC) methods are a class of algorithms for *sampling from a probability distribution*.

They are **based on a Markov chain** that has the desired distribution of its **equilibrium** distribution.

MCMC algorithms are generally used for sampling from multi-dimensional distributions, especially when the *number of dimensions is high*.

Our roadmap today:

- 1. What is a Markov chain?
- 2. How do we sample from a Markov chain?
- 3. How to use MCMC for sampling from target distribution 1
- 4. How to use MCMC for sampling from target distribution 2

Unit 3 Markov Chain

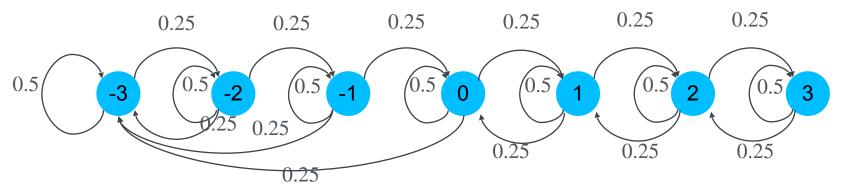
Unit 4 MCMC

Unit 5 Gibbs Samplir
Unit 6 Metropolis

Hastings

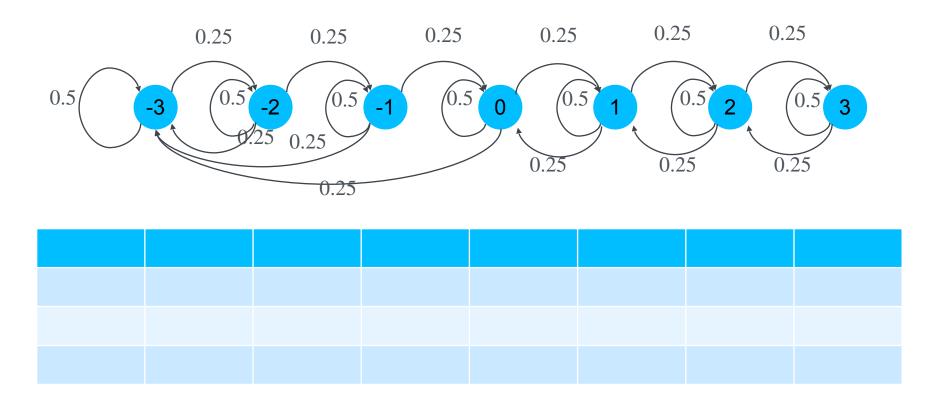
Markov Chain

Example



Definition: A Markov chain defines a probabilistic transition model $T(s \rightarrow s')$ over states S:

For all
$$s \in S$$
:
$$\sum_{s'} T(s \to s') = 1$$



Stationary Process

General definition: A stochastic process $\left\{X^{(t)}: t \geq 0\right\}$ is stationary if for any time points i_1, \ldots, i_n and any $m \geq 0$, the joint distribution of $(X^{(i_1)}, \ldots, X^{(i_n)})$ is the same as the joint distribution of the time-shifted $(X^{(m+i_1)}, \ldots, X^{(m+i_n)})$.

Implies: If a process $\{X^{(t)}: t \geq 0\}$ is stationary, the distributions $X^{(t)}$ are the same for all t.

Stationary Distribution π

$$P^{(t)}(s') \approx P^{(t+1)}(s') = \sum_{s \in S} P^{(t)} (X^{(t)} = s) T(s \to s')$$

$$\pi(s') = \sum_{s \in S} \pi(s) T(s \to s')$$

• Matrix $T \in \mathbb{R}^{|S| \times |S|}$ is a stochastic matrix:

$$\forall s \in S: \quad \sum_{s' \in S} T[s \to s'] = 1$$

• $\pi(s)$ is constant for all $s \in S$,

this leads to the following eigenvalue problem:

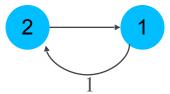
$$\pi^{\dagger}\lambda = \pi^{\dagger} T$$
, with $\lambda = 1$

Regular Markov Chain

Definition: A Markov chain is regular if there exists k such that, for every s, s' the probability of getting from s to s' in exactly k steps is larger than 0.

Counterexample:

1



Regular Markov Chain

Definition: A Markov chain is regular if there exists k such that, for every s, s' the probability of getting from s to s' in exactly k steps is larger than 0.

Theorem: A regular Markov chain converges to a unique stationary distribution regardless of start state.

Making a Markov Chain regular

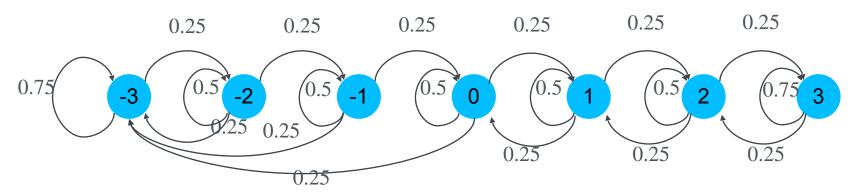
If one can transition from every state in a Markov Chain to every other state (and it is not regular yet), you can make it regular by adding self transitions.

$$T^{\text{new}} = \frac{1}{2}(T+I)$$



4 Markov Chain Monte Carlo

Sampling from a Regular Markov Chain



Very simple idea:

- 1. Start in a randomly chosen state
- 2. Repeat until probability distribution is stationary
 - Randomly transition into next state s' with probability $T[s \longrightarrow s']$
- 3. Repeat until sample is large enough
 - sample current state
 - Randomly transition into next state s' with probability $T[s \longrightarrow s']$

Using a Markov Chain for a Non-Markov Chain distribution

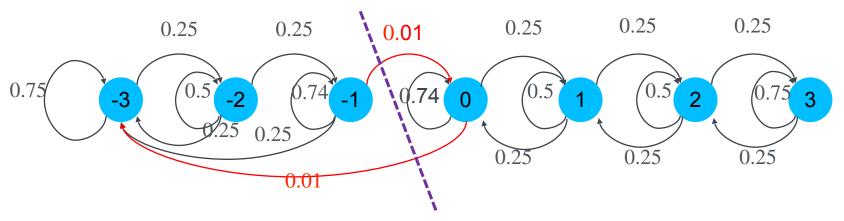
- Goal: compute $P(x \in X)$
 - but: P(X) is too hard to sample from directly
- Approach:
 - 1. Construct a Markov chain T whose unique stationary distribution is P(X)
 - 2. Sample $x^{(0)}$ from some $P^{(0)}$
 - 3. For t = 0, 1, 2, ...
 - Generate $x^{(t+1)}$ from $T(x^{(t)} \longrightarrow x')$

- Considerations
 - Only sample if distribution is close to ${m P}({m X})$
 - At early iterations, $P^t(X)$ is usually far from $P(X) = \pi(X)$
 - Thus: collect samples only after the chain has run long enough to *mix*

Mixing

Common mixing problem

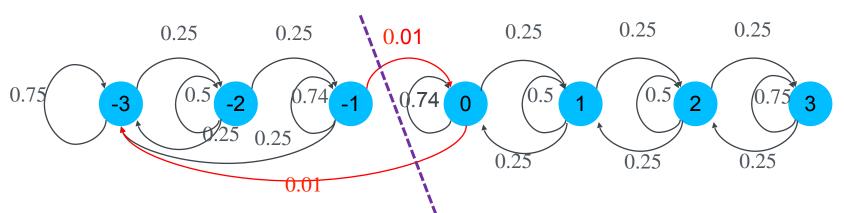
samples are not independent!



- How do you know if a chain has mixed or not?
 - In general, you cannot prove a chain has mixed
 - Often, you can show that is has NOT mixed

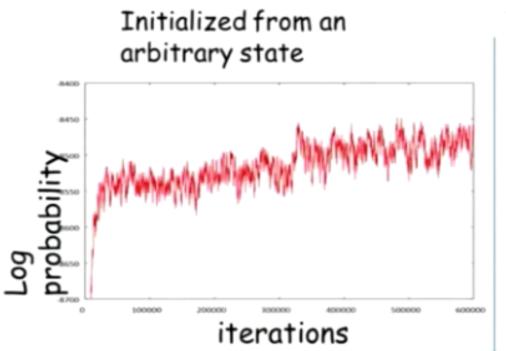
Mixing

Common mixing problem

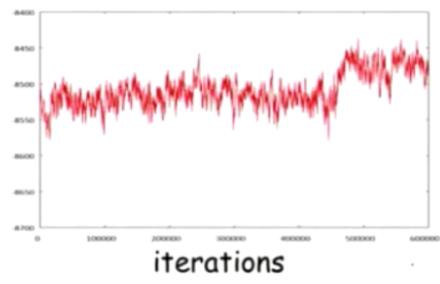


- How do you know a chain has not mixed?
 - Compare chain statistics in different windows within a single run of the chain
 - and across different runs, initialized differently

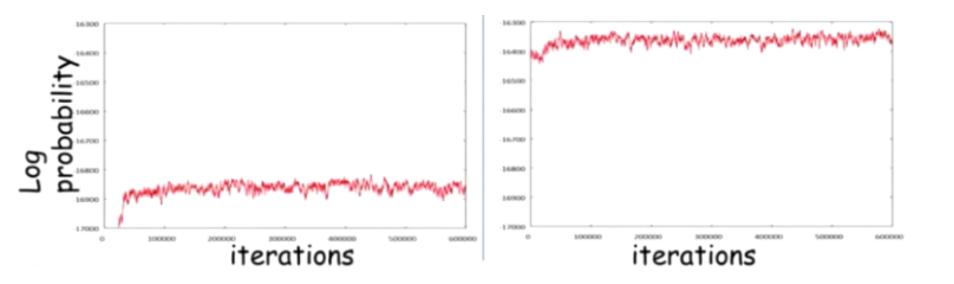
Comparing two runs



Initialized from a highprobability state



Comparing two runs



Tools

- All kind of statistics
 - mean
 - variance
- All kind of visualizations
 - scatter plots mapping two series

• ...

About the Samples

- Once the chain mixes (at $t^{(\text{mix})}$), all samples $x^{(t)}$ are from the stationary distribution π
 - Use all $x^{(t)}$ for $t > t^{(\text{mix})}$
- But: nearby samples are correlated (not iid!)
 - do not overestimate quality of sample, just based on size
- When a chain mixes faster,
 the samples are less correlated and more useful

MCMC (refined)

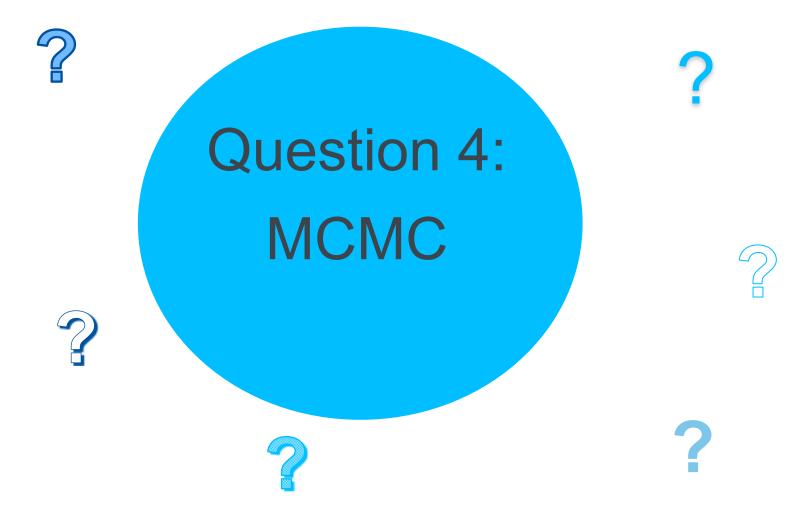
- 1. Iterate over runs r = 1...R
 - Sample $x^{(r,0)}$ from $P^{(0)}$
- 2. Repeat until mixing
 - For $r = 1 \dots R$
 - Generate $x^{(r,t+1)}$ from $T(x^{(r,t)} \rightarrow x')$
 - Compare window statistics in different chains to determine mixing

- 3. Repeat until sufficient samples
 - $D \coloneqq \emptyset$
 - For r = 1...R
 - Generate $x^{(r,t+1)}$ from $T(x^{(r,t)} \rightarrow$
 - $D \coloneqq D \cup \left\{ x^{(r,t+1)} \right\}$
 - t := t + 1
- $4. \quad D = \{x[1] \dots x[M]\}$

$$\mathbb{E}_{P(\mathcal{X})}[f] \approx \frac{1}{M} \sum_{m=1}^{M} f(x[m])$$

Issues with MCMC

- Many parameters
 - how many samples to collect?
 - how correlated are they?
 - how long until mixing?
 - how to compute statistics that determines mixing?
 - window size
 - which test statistics?
 - how close should test statistics be?



5 Gibbs Sampling

How to design a markov chain? Answer 1: Gibbs chain

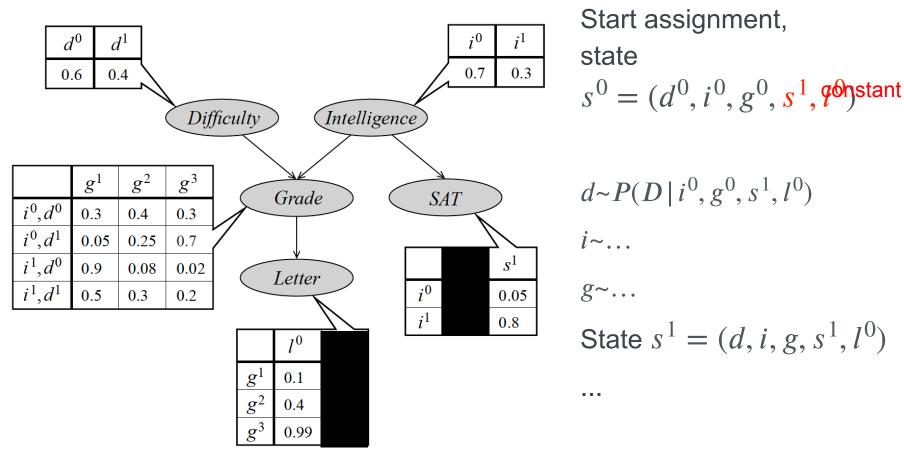
- Target distribution $\mathbf{\textit{P}}_{\Phi}(\mathcal{X}) = \mathbf{\textit{P}}_{\Phi}(X_1,...,X_n)$
- Markov chain state space: each complete assignment x to $\mathcal X$ is a state of the Markov chain
- Let x_{-i} denote $(x_1...x_{i-1} x_{i+1} ...x_n)$

- ullet Random transition step defined given state $oldsymbol{x}$
 - For i = 1...n
 - Sample $x_i' \sim P_{\Phi}(X_i \mid x_{-i})$
 - Set $x = x_{-i} \circ x_i'$
 - Return x'=x

No explicit definition of T

Gibbs sampling with evidence

Sample $P_{\Phi}(D, I, G | s^1, l^0)$



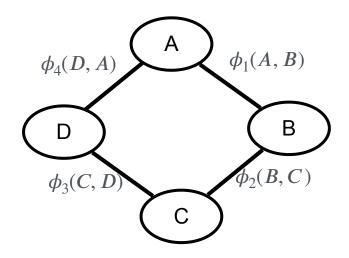
Computational cost

$$P_{\Phi}(A \mid b, c, d) = \frac{\tilde{P}_{\Phi}(A, b, c, d)}{\sum_{A'} \tilde{P}_{\Phi}(A', b, c, d)} = \frac{\phi_1(A, b)\phi_2(b, c)\phi_3(c, d)\phi_4(d, A)}{\sum_{A'} \phi_1(A', b)\phi_2(b, c)\phi_3(c, d)\phi_4(d, A')} = \frac{\phi_1(A, b)\phi_4(d, A)}{\sum_{A'} \phi_1(A', b)\phi_4(d, A')} = \frac{\phi_1(A, b)\phi_4(d, A')}{\sum_{A'} \phi_1(A', b)\phi_4(d, A')} = \frac{\phi_1(A, b)\phi_4(A', b)\phi_4(d, A')}{\sum_{A'} \phi_1(A', b)\phi_4(A', b)} = \frac{\phi_1(A, b)\phi_4(A', b)\phi_4(A', b)}{\sum_{A'} \phi_1(A', b)\phi_4(A', b)} = \frac{\phi_1(A, b)\phi_4(A', b)\phi_4(A', b)}{\sum_{A'} \phi_1(A', b)\phi_4(A', b)} = \frac{\phi_1(A', b)\phi_4(A', b)\phi_4(A', b)}{\sum_{A'} \phi_1(A', b)\phi_4(A', b)} = \frac{\phi_1(A', b)\phi_4(A', b)}{\sum_{A'} \phi_1(A', b)\phi_4(A', b)} = \frac{\phi$$

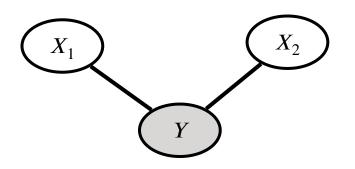
- For i = 1...n
 - Sample $x_i' \sim P_{\Phi}(X_i | x_{-i})$

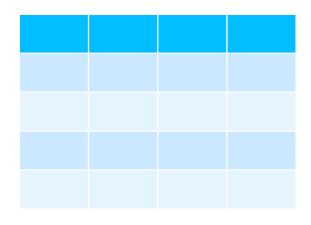
$$P_{\Phi}(X_i \mid x_{-i}) = \frac{P_{\Phi}(X_i, x_{-i})}{P_{\Phi}(x_{-i})} = \frac{\frac{1}{Z} \tilde{P}_{\Phi}(X_i, x_{-i})}{\frac{1}{Z} \tilde{P}_{\Phi}(x_{-i})} = \frac{\tilde{P}_{\Phi}(X_i, x_{-i})}{\tilde{P}_{\Phi}(x_{-i})}$$

Only have to consider factors related to X_i



Counterexample to regularity: XOR



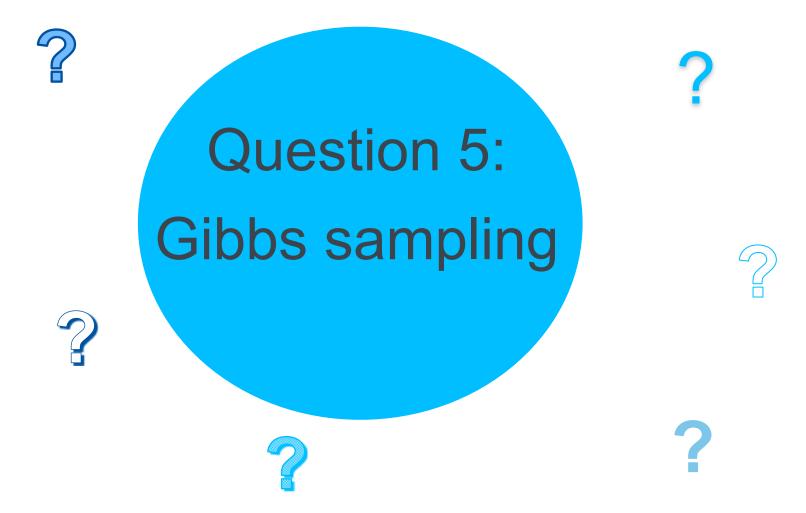


• Compute $P(X_1, X_2 | Y = 1)$

If all factors are positive, Gibbs chain is regular (connected, aperiodic)

Conclusion of Gibbs sampling

- Handles high-dimensional problem
- Handles continuous and discrete random variables
- Sampling from conditionals can be hard!
 - Possible way out: "Metropolis within Gibbs"
- If random variables are strongly correlated, the Markov chain converges slowly and has highly auto-correlated samples



6 Metropolis Hastings Algorithm

How to apply MCMC?

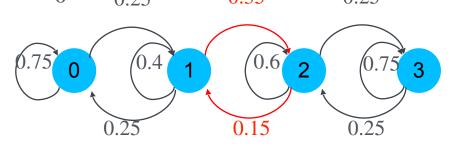
- Gibbs:
 - states correspond to complete assignments
 - sample from posteriors $P(X_i | \mathbf{x}_{-i})$

- Metropolis Hastings:
 - sample proportional to $extbf{\emph{P}}_{m{\Phi}}(m{\mathscr{X}})$
 - ullet unnormalized $oldsymbol{\widetilde{P_{\Phi}}}(\mathcal{X})$!
 - can also be proportional to density for continuous probability distributions

Different strengths and weaknesses wrt mixing.
No single best choice!

Reversible Markov Chain

Definition: A Markov chain is called reversible, if for every n the sequence $X_0, X_1, ... X_n$ has the same joint probability as the time-reversed sequence $X_n, X_{n-1}, ... X_0$.



In particular this implies:

Theorem: A Markov chain with probability distribution over states π is reversible if and only if the detailed balance equation holds:

$$\pi(x)T(x \longrightarrow x') = \pi(x')T(x' \longrightarrow x)$$

Reversible Markov Chain

$$\pi(x)T(x \longrightarrow x') = \pi(x')T(x' \longrightarrow x)$$

Theorem: If T is regular, it has a unique stationary distribution.

Theorem: If T is regular and detailed balance holds, the unique stationary distribution is π .

Proof:

$$\sum_{x} \pi(x)T(x \longrightarrow x') = \sum_{x} \pi(x')T(x' \longrightarrow x)$$

$$\sum_{x} \pi(x)T(x \longrightarrow x') = \pi(x')\sum_{x} T(x' \longrightarrow x) = \pi(x')$$

Metropolis Hastings Chain

- Proposal distribution $Q(x \to x') \in [0,1]$
 - describes how to explore the space ${\mathscr X}$
- Acceptance probability $A(x \to x')$ Bernoulli random variable
 - reflects $P_{oldsymbol{\Phi}}({\mathscr X})$
 - ullet Turns ergodic Markov chain Q into a reversible Markov chain T
- Targeted Markov chain

$$T(x \to x') = Q(x \to x')A(x \to x'), \text{ if } x \neq x'$$

$$T(x \to x') = Q(x \to x') + \sum_{x' \in X} Q(x \to x')(1 - A(x \to x')), \text{ else}$$

ullet Meaning: Q proposes transition to next sample and A accepts it or not

Make Metropolis Hastings Chain Reversible

- Proposal distribution $Q(x \to x') \in [0,1]$
- Acceptance probability $A(x \to x') \in \{0,1\}$
- Targeted Markov chain: $T(x \to x') = Q(x \to x')A(x \to x')$, if $x \neq x'$
- Enforce reversibility

$$\pi(x)T(x \longrightarrow x') = \pi(x')T(x' \longrightarrow x)$$

$$\pi(x)Q(x \to x')A(x \to x') = \pi(x')Q(x' \to x)A(x' \to x)$$

$$\frac{A(x \to x')}{A(x' \to x)} = \frac{\pi(x')Q(x' \to x)}{\pi(x)Q(x \to x')} = p < 1, \text{ set } A(x' \to x) = 1$$

$$A(x \to x') = \min\left(1, \frac{\pi(x')Q(x' \to x)}{\pi(x)Q(x \to x')}\right)$$

Deliver $P_{oldsymbol{\Phi}}({oldsymbol{\mathscr{X}}})$

- Proposal distribution $Q(x \to x') \in [0,1]$
- Acceptance probability $A(x \to x') \in \{0,1\}$
- Targeted Markov chain:

$$T(x \to x') = Q(x \to x')A(x \to x')$$
, if $x \neq x'$

Enforce reversibility

$$A(x \to x') = \min\left(1, \frac{\pi(x')Q(x' \to x)}{\pi(x)Q(x \to x')}\right)$$

$$\frac{\pi(x')}{\pi(x)} = \frac{P_{\Phi}(x')}{P_{\Phi}(x)}$$

Choice of Q

$$A(x \to x') = \min\left(1, \frac{P_{\Phi}(x')Q(x' \to x)}{P_{\Phi}(x)Q(x \to x')}\right)$$

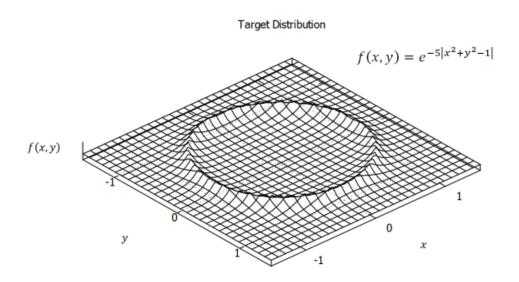
• *Q* must be reversible:

•
$$Q(x' \to x) > 0 \Longrightarrow Q(x \to x') > 0$$

- ullet describes how to explore the space ${\mathscr X}$
 - movements should be neither to small nor too large
 - large "jumps" often end in area of low probability

 high rejection rate
 - ullet often a normal distribution (for continuous distributions ${\mathscr X}$)

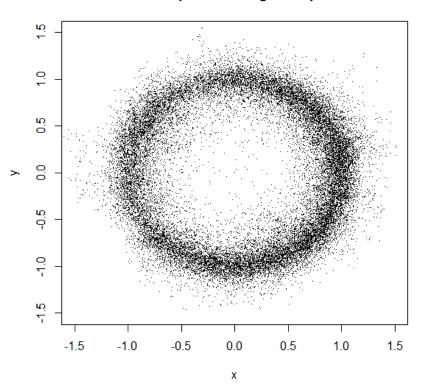
Illustration Metropolis-Hastings in Continuous Space



From: https://blog.abhranil.net/2014/02/08/r-code-for-multivariate-random-walk-metropolis-hastings-sampling/

Illustration Metropolis-Hastings

Metropolis-Hastings Sample



From: https://blog.abhranil.net/2014/02/08/r-code-for-multivariate-random-walk-metropolis-hastings-sampling/

Lessons learned - MH Sampling is susceptible to:

- Start points
 - AlphaGo uses neural network just to predict good start point for sampling
- Proposal distribution
 - · Proposal steps should not be too large or too small
 - Too large many rejections
 - Too small slow diffusion
 - Tremendous flexibility in designing proposal distributions that explore the space quickly
 - · good proposal distribution makes a huge difference
- Mixing time (also called burn-in time)
 - Do not sample too early
- Lag (also referred to as sampling intervals)
 - · Succeeding states are closely coupled
- Shape of target function
 - Non-differentiable is problematic
 - Low or zero slopes are problematic
- Discrete distributions can be harder to sample from than continuous ones

Metropolis Hastings Algorithm

Also used to compute

- expected values
- integrals

Metropolis Hastings is a general framework for building Markov chains with a particular stationary distribution

First version developed in 1950 in statistical physics





Thank you!



Steffen Staab

E-Mail Steffen.staab@ipvs.uni-stuttgart.de
Telefon +49 (0) 711 685e be defined
www.ipvs.uni-stuttgart.de/departments/ac/

Universität Stuttgart Analytic Computing, IPVS Universitätsstraße 32, 50569 Stuttgart