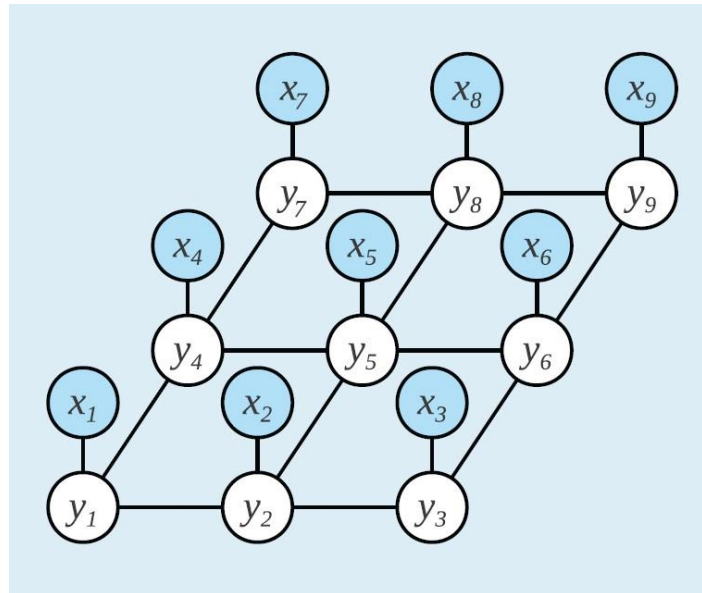


Probabilistic Graphical Models in Bioinformatics

Lecture 15: Particle-based approximate inference



Introduction

- Particle-based methods approximate the joint distribution as a set of random samples from the network.
- The term “particle” came from statistical physics.
- Outline
 - Forward sampling
 - Likelihood weighting and importance sampling
 - Markov Chain Monte Carlo

Forward sampling

- We generate random samples $\xi[1], \dots, \xi[M]$ from the distribution $P(\mathcal{X})$
- We then use the samples to compute the expectation of some target function f .

$$\hat{E}_{\mathcal{D}}(f) = \frac{1}{M} \sum_{m=1}^M f(\xi[m]).$$

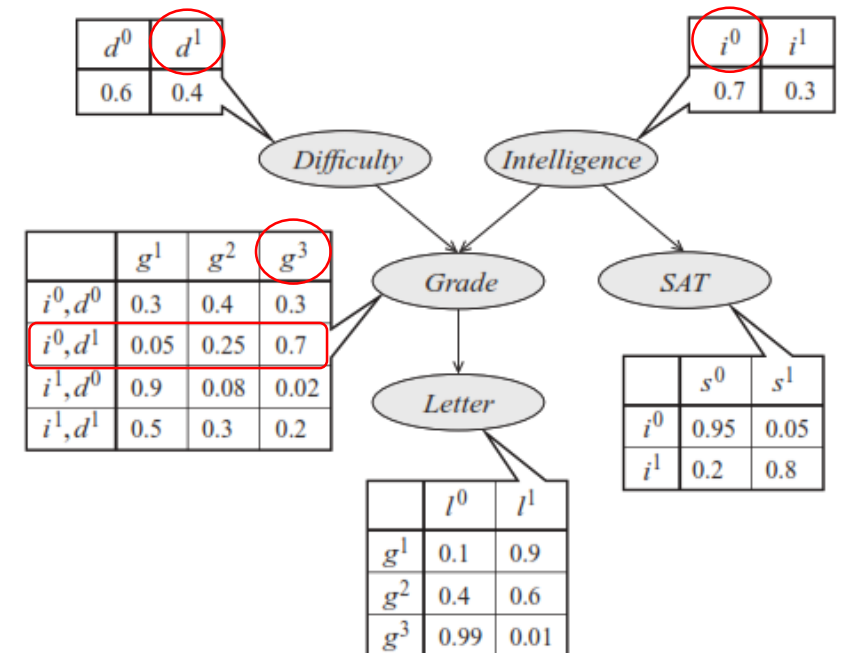
$$\hat{P}_{\mathcal{D}}(\mathbf{y}) = \frac{1}{M} \sum_{m=1}^M \mathbf{1}\{\mathbf{y}[m] = \mathbf{y}\},$$

$\mathbf{y}[m]$ denotes assignment to variables \mathbf{Y} in the particles $\xi[m]$

Or simply the fraction of particles with the event \mathbf{y} observed.

Forward sampling- sampling from a Bayesian network

- A very simple process
- We sample the nodes in some order consistent with the partial order of the BN
- So by the time we sample a node we have values for all of its parents.
- We can then sample from the distribution defined by the CPD and by the chosen values for node's parents.



Analysis of error

- The quality of the estimate obtained depends on the number of particles generated.
- *The Hoeffding bound*
 - guarantees a bound on the absolute error

$$P_{\mathcal{D}}(\hat{P}_{\mathcal{D}}(\mathbf{y}) \notin [P(\mathbf{y}) - \epsilon, P(\mathbf{y}) + \epsilon]) \leq 2e^{-2M\epsilon^2}.$$

- The number of samples need to bound error by ϵ , with probability at least $1 - \delta$.

$$M \geq \frac{\ln(2/\delta)}{2\epsilon^2}.$$

- The *Chernoff bound*
 - guarantees a bound on the relative error

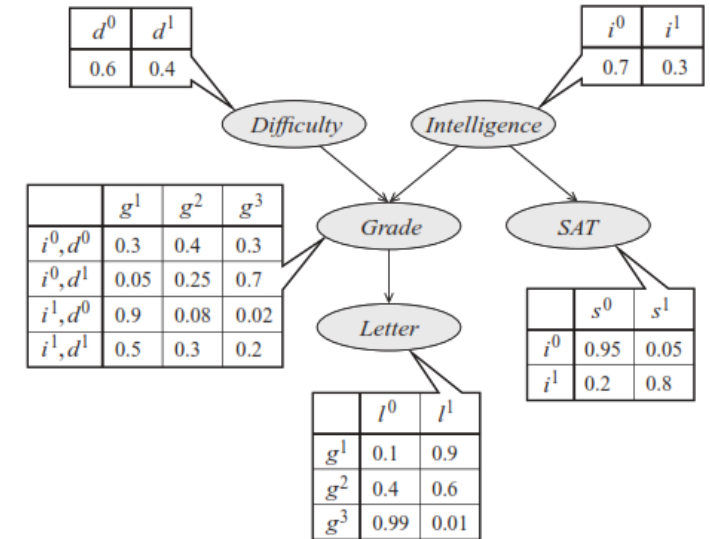
$$P_{\mathcal{D}}(\hat{P}_{\mathcal{D}}(\mathbf{y}) \notin P(\mathbf{y})(1 \pm \epsilon)) \leq 2e^{-MP(\mathbf{y})\epsilon^2/3}.$$

Conditional probability queries

- So far we discuss estimating marginal probabilities
- In general we are interested in conditional probabilities of the form $P(y \mid E = e)$.
 - Much harder
- Rejection sampling
 - we generate samples x from the $P(\mathcal{X})$
 - and then reject any sample that is not compatible with e
 - the resulting samples are from posterior $P(\mathcal{X} \mid e)$
- Issue with rejection sampling
 - The number of unrejected samples can be quite small (i.e., $MP(e)$).
 - **Example:** if $P(e) = 0.001$ we need to sample $M = 10,000$ to on average obtain 10 unrejected particles.
- Low probability of evidence is the rule than the exception.

Forcing the samples to match evidence

- Assume that our evidence is d^1, s^1
 - forward sampling process might generate a value of d^0 for D
 - this sample will always be rejected as being incompatible with the evidence.
- A more sensible approach is
 - to simply force the samples to match evidence
 - without correction, this approach can generate incorrect results.
- Example
 - Assume the evidence is s^1
 - Using the naïve process, we first sample D and I , set $S = s^1$, and then sample G and L appropriately.
 - What is the expected fraction of samples that have $I = i^1$?
 - How do you compare it with $P(I = i^1 \mid s^1)$?



Likelihood weighting

- Likelihood weighting algorithm generated weighted particles according to the likelihood of the evidence (*see section 12.2.1 for some examples and intuition*)

$$\langle \xi[1], w[1] \rangle, \dots, \langle \xi[M], w[M] \rangle.$$

- We then estimate

$$\hat{P}_{\mathcal{D}}(\mathbf{y} \mid \mathbf{e}) = \frac{\sum_{m=1}^M w[m] \mathbf{I}\{\mathbf{y}[m] = \mathbf{y}\}}{\sum_{m=1}^M w[m]}.$$

Algorithm 12.2 Likelihood-weighted particle generation

```

Procedure LW-Sample (
     $\mathcal{B}$ , // Bayesian network over  $\mathcal{X}$ 
     $\mathbf{Z} = \mathbf{z}$  // Event in the network
)
1  Let  $X_1, \dots, X_n$  be a topological ordering of  $\mathcal{X}$ 
2   $w \leftarrow 1$ 
3  for  $i = 1, \dots, n$ 
4       $\mathbf{u}_i \leftarrow \mathbf{x} \langle \text{Pa}_{X_i} \rangle$  // Assignment to  $\text{Pa}_{X_i}$  in  $x_1, \dots, x_{i-1}$ 
5      if  $X_i \notin \mathbf{Z}$  then
6          Sample  $x_i$  from  $P(X_i \mid \mathbf{u}_i)$ 
7      else
8           $x_i \leftarrow \mathbf{z} \langle X_i \rangle$  // Assignment to  $X_i$  in  $\mathbf{z}$ 
9           $w \leftarrow w \cdot P(x_i \mid \mathbf{u}_i)$  // Multiply weight by probability of desired value
10 return  $(x_1, \dots, x_n), w$ 

```

Importance sampling

- Importance sampling is a general approach for estimating the expectation of a function $f(x)$ relative to some distribution $P(X)$, typically called the *target distribution*.
- We can obtain estimates of this expectation by generating samples from a different distribution Q .

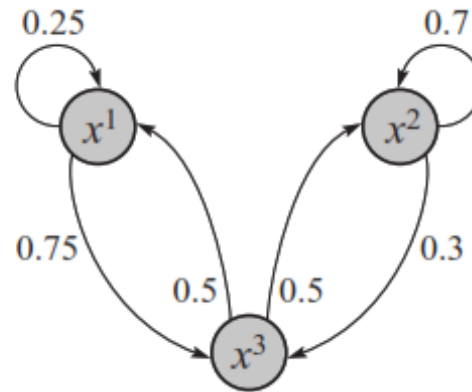
$$\hat{E}_{\mathcal{D}}(f) = \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}[m]) \frac{P(\mathbf{x}[m])}{Q(\mathbf{x}[m])}.$$

- Reasons to sample from a different distribution:
 - It might be impossible or computationally very expensive to generate samples from P .

Markov Chain Monte Carlo Methods

Markov chain

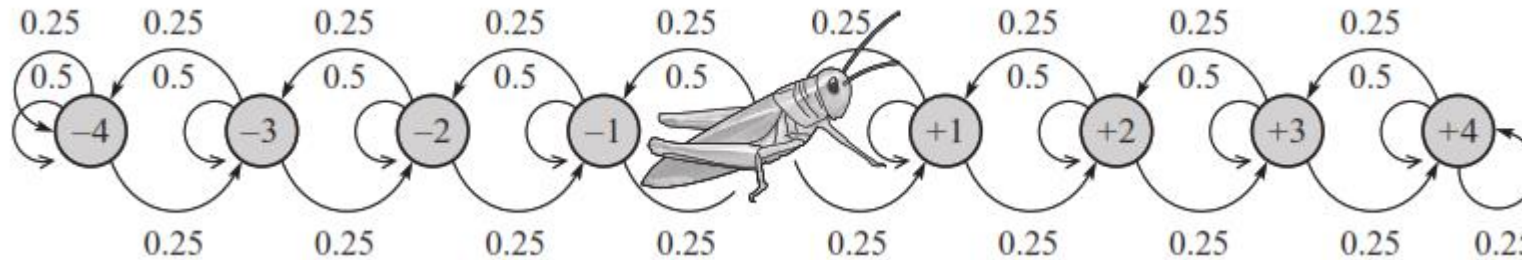
- A Markov chain is defined in terms of a graph of states over which the sampling algorithm takes a random walk
- It defines a probabilistic transition model $\mathcal{T}(x \rightarrow x')$ over states x .



- Random sampling process
 - Defines a random sequence of states $x^{(0)}, x^{(1)}, x^{(2)}, \dots$
 - Because the transition model is random, the state of the process at step t can be viewed as a random variable $X^{(t)}$
 - Assume the initial state $X^{(0)}$ is distributed according to some initial state distribution $P^{(0)}(X^{(0)})$

$$P^{(t+1)}(X^{(t+1)} = x') = \sum_{x \in \text{Val}(X)} P^{(t)}(X^{(t)} = x) \mathcal{T}(x \rightarrow x').$$

Example: Grasshopper Markov chain



	-2	-1	0	+1	+2
$P^{(0)}$	0	0	1	0	0
$P^{(1)}$	0	0.25	0.5	0.25	0
$P^{(2)}$			$.5^2 + 2 \times 0.25^2$ $= 0.375$		
$P^{(50)}$	At t=50, the distribution is almost uniform with a range of 0.1107-0.1116				

The most important aspect of a Markov chain is its long-term behavior

Markov chain Monte carlo (MCMC) sampling



- MCMC sampling is a process that mirrors the dynamics of the Markov chain

Algorithm 12.5 Generating a Markov chain trajectory

```
Procedure MCMC-Sample (  
     $P^{(0)}(\mathbf{X})$ , // Initial state distribution  
     $\mathcal{T}$ , // Markov chain transition model  
     $T$  // Number of time steps  
)  
1  Sample  $\mathbf{x}^{(0)}$  from  $P^{(0)}(\mathbf{X})$   
2  for  $t = 1, \dots, T$   
3      Sample  $\mathbf{x}^{(t)}$  from  $\mathcal{T}(\mathbf{x}^{(t-1)} \rightarrow \mathbf{X})$   
4      return  $\mathbf{x}^{(0)}, \dots, \mathbf{x}^{(T)}$ 
```

- The sample $\mathbf{x}^{(t)}$ is drawn from the distribution $P^{(t)}$.
- We are interested in the limit of this process
 - whether $P^{(t)}$ converges
 - and if so, to what limit

Stationary distribution

- Intuitively as the process converges, we would expect $P^{(t+1)}$ to be close to $P^{(t)}$

$$P^{(t)}(\mathbf{x}') \approx P^{(t+1)}(\mathbf{x}') = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} P^{(t)}(\mathbf{x}) \mathcal{T}(\mathbf{x} \rightarrow \mathbf{x}').$$

- A distribution $\pi(\mathbf{X})$ is a stationary distribution for a *Markov chain* if it satisfies

$$\pi(\mathbf{X} = \mathbf{x}') = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \pi(\mathbf{X} = \mathbf{x}) \mathcal{T}(\mathbf{x} \rightarrow \mathbf{x}').$$

- We already discussed the uniform distribution is a stationary distribution for the Grasshopper Markov chain.
- In general there is no guarantee the stationary distribution is unique.

Stationary distribution-example

- By definition we have the following equations:

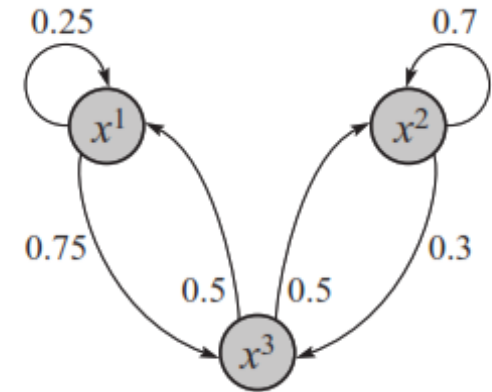
$$\begin{aligned}\pi(x^1) &= 0.25\pi(x^1) + 0.5\pi(x^3) \\ \pi(x^2) &= 0.7\pi(x^2) + 0.5\pi(x^3) \\ \pi(x^3) &= 0.75\pi(x^1) + 0.3\pi(x^2),\end{aligned}$$

- in addition, it needs to be a legal distribution:

$$\pi(x^1) + \pi(x^2) + \pi(x^3) = 1.$$

- it is easy to verify the system has a unique solution

$$\begin{aligned}\pi(x^1) &= 0.2 \\ \pi(x^2) &= 0.5 \\ \pi(x^3) &= 0.3\end{aligned}$$



Regular Markov chains

- A Markov chain is said to be regular if there exists some number k such that the probability of getting between every two states in exactly k steps is > 0 .
- Theorem: if a finite state Markov chain is **regular**, then it has a unique stationary distribution.
- Sufficient conditions for regularity
 - There is a positive probability path between every two states in the state graph.
 - For every state, there is a self-loop.

Using a Markov chain for answering
probabilities queries

Using a Markov chain

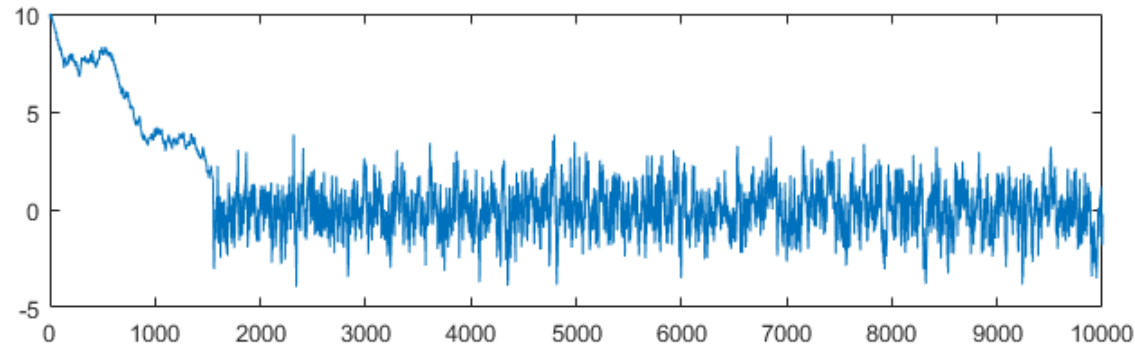
- The theory of Markov chains provides a general framework for generating samples from a target distribution π .
- Goal: compute $P(X | E = e)$
 - Too difficult to directly sample from it.
- To this end, we need to build a Markov chain whose unique stationary distribution is $P(X | E = e)$
 - We need to define an appropriate transition model

```
Sample  $x^{(0)}$  from  $P^{(0)}(X)$ 
for  $t = 1, \dots, T$ 
  Sample  $x^{(t)}$  from  $\mathcal{T}(x^{(t-1)} \rightarrow X)$ 
```

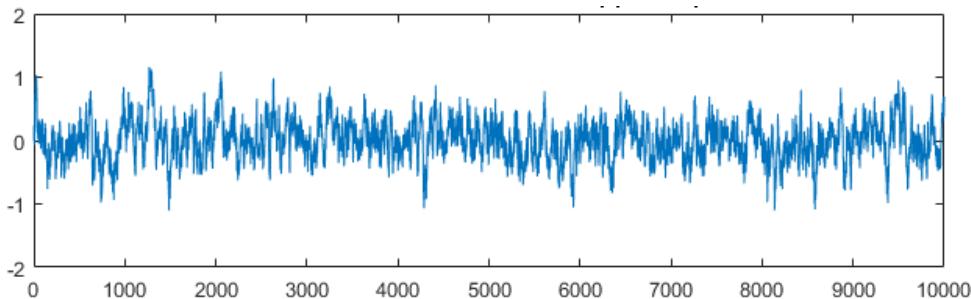
- At early step, $P^{(t)}$ is very different from the P , the stationary distribution

Mixing

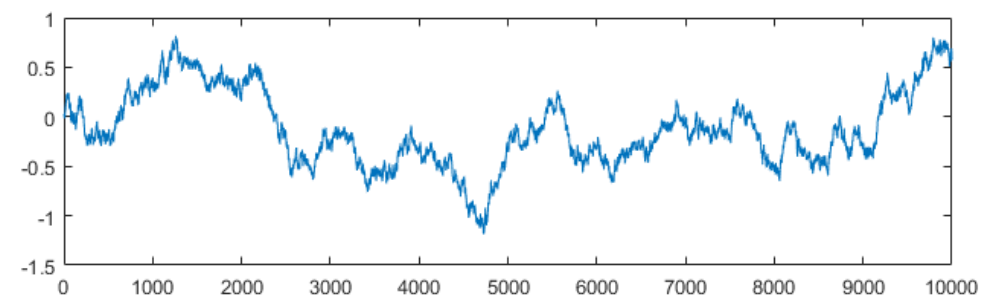
- **The mixing time** of a Markov chain is the time until the Markov chain is "close" to its stationary state distribution.
- *Burn-in time*: the number of steps we take until we collect a sample from the chain.



- How do we evaluate the time required for the a chain to mix?
 - There is no general-purpose theoretical analysis for the mixing time. However,
 - we can compare chain statistics in different windows of a single run
 - or across different runs of the chain



Well-mixed



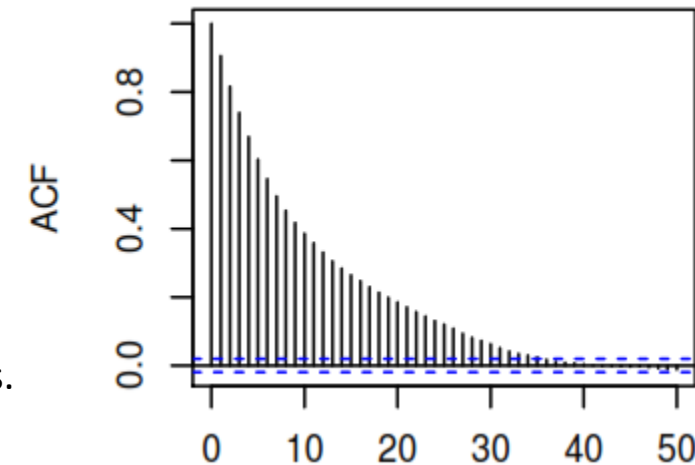
Highly correlated samples

Using the samples

- After discarding the samples in the burn-in phase, we collect M samples from the chain.
- Then we can estimate expectation of any function f by

$$\hat{E}_D(f) = \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}[m], e)$$

- $\hat{E}_D(f)$ is an unbiased estimator $\hat{E}_{\pi(X)}[f(X, e)]$.
- The samples in the chain are correlated
 - If we are looking for a set of independent samples, we only take every n th sample
 - Variable n is often determined by examining autocorrelation between adjacent samples.
- In many settings, there is no need to look for an independent set of samples
 - Even though correlated, using all samples produces a better estimator



Metropolis-Hastings algorithm

Reversible Markov chains

- Detailed balance equation:

$$\pi(\mathbf{x})\mathcal{T}(\mathbf{x} \rightarrow \mathbf{x}') = \pi(\mathbf{x}')\mathcal{T}(\mathbf{x}' \rightarrow \mathbf{x}).$$

- **Theorem:** if \mathcal{T} is regular and it satisfies the detailed balance equation relative to π , then π is the unique stationary distribution of \mathcal{T} .
- The reversibility condition gives us a condition for verifying that our Markov chain has the desired stationary distribution.
- Metropolis-Hastings algorithm is a general method to build a reversible Markov chain with a particular stationary distribution.
 - Uses the idea of a proposal distribution

Metropolis-Hastings algorithm

- Goal: to design a reversible Markov chain for sampling $\pi(x)$
 - Sampling directly from $\pi(x)$ is hard but it is possible to evaluate it.
- Proposal distribution $T^Q(x \rightarrow x')$
- Acceptance probability $A(x \rightarrow x')$
- At each state x , sample x' from $T^Q(x \rightarrow x')$
- Accept proposal with probability $A(x \rightarrow x')$
 - If proposal accepted, move to x' .
 - Otherwise stay at x
- The actual transition model

$$\begin{aligned}\mathcal{T}(x \rightarrow x') &= \mathcal{T}^Q(x \rightarrow x')\mathcal{A}(x \rightarrow x') & x \neq x' \\ \mathcal{T}(x \rightarrow x) &= \mathcal{T}^Q(x \rightarrow x) + \sum_{x' \neq x} \mathcal{T}^Q(x \rightarrow x')(1 - \mathcal{A}(x \rightarrow x')).\end{aligned}$$

Metropolis-Hastings algorithm-2

- Given a proposal distribution, we can use the detailed balance to find $A(x \rightarrow x')$ such that the stationary distribution is $\pi(x)$
- The detailed balance equation assert

$$\pi(x)\mathcal{T}^Q(x \rightarrow x')\mathcal{A}(x \rightarrow x') = \pi(x')\mathcal{T}^Q(x' \rightarrow x)\mathcal{A}(x' \rightarrow x).$$

- Hence we obtain

$$\mathcal{A}(x \rightarrow x') = \min \left[1, \frac{\pi(x')\mathcal{T}^Q(x' \rightarrow x)}{\pi(x)\mathcal{T}^Q(x \rightarrow x')} \right],$$