

Dicing

Chen Yunchuan
2015.4.29

Outline

- Sampling
 - Random Variable Generation
 - Rejection Sampling
 - Importance Sampling
 - Markov Chain Monte Carlo (MCMC)
- Restricted Boltzmann Machine Revisit
 - Learning algorithms
- Imagination: Represent Word with Random Vectors

Example of Monte Carlo: Estimate π

```
procedure direct-pi
   $N_{\text{hits}} \leftarrow 0$  (initialize)
  for  $i = 1, \dots, N$  do
    {
       $x \leftarrow \text{ran}(-1, 1)$ 
       $y \leftarrow \text{ran}(-1, 1)$ 
      if  $(x^2 + y^2 < 1)$   $N_{\text{hits}} \leftarrow N_{\text{hits}} + 1$ 
    }
  output  $N_{\text{hits}}$ 
```



- Estimate π

$$\pi \approx \frac{4N_{\text{hits}}}{N}$$

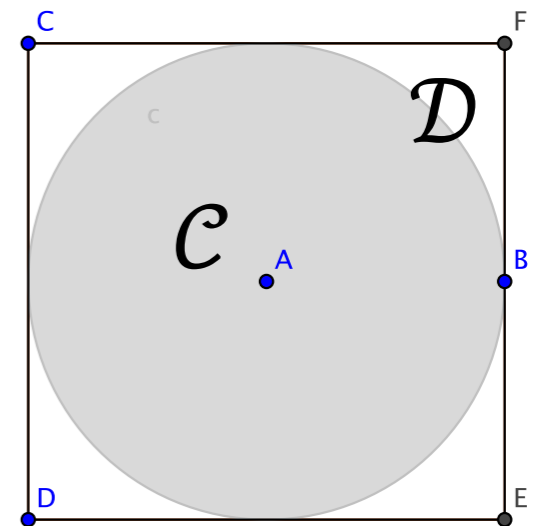
Run	N_{hits}	Estimate of π
1	3156	3.156
2	3150	3.150
3	3127	3.127
4	3171	3.171
5	3148	3.148

Estimating π

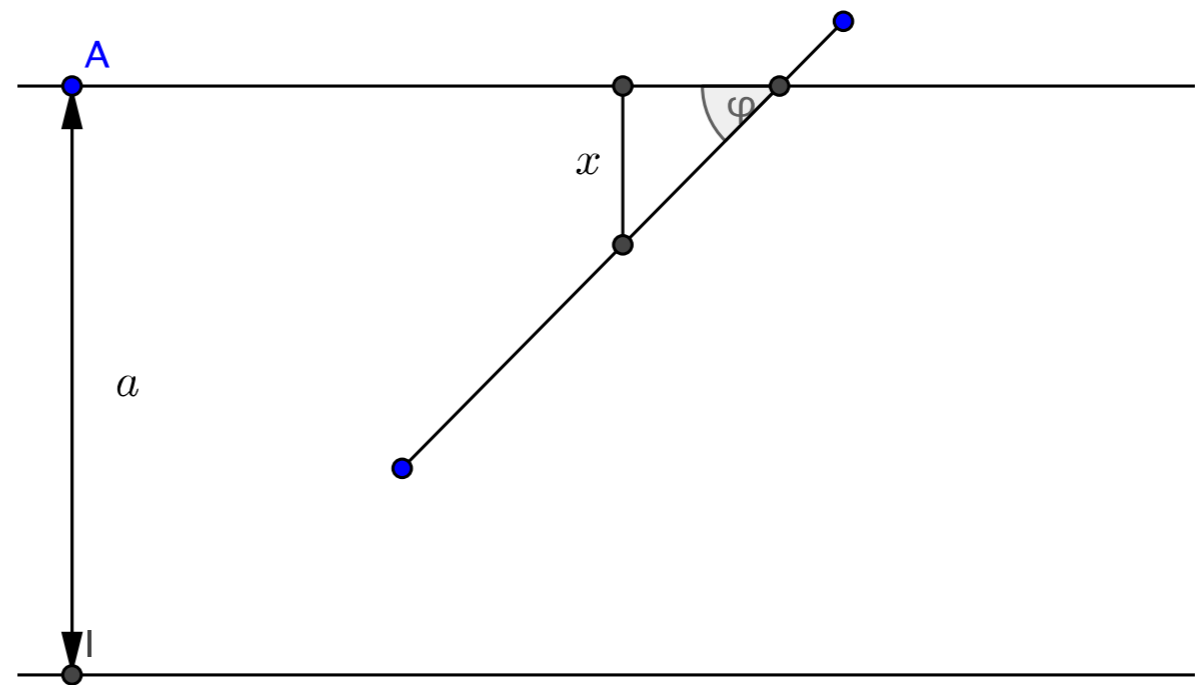
Draw i.i.d. set of samples $\{\mathbf{x}^{(i)}\}_{i=1}^N$
from a target density $p(\mathbf{x}|\cdot)$

$$\int_{\mathcal{X}} f(\mathbf{x})p(\mathbf{x}|\cdot)d\mathbf{x} \approx \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}^{(i)})$$

$$\begin{aligned} \frac{\pi}{4} &= \frac{S_{\mathcal{C}}}{S_{\mathcal{D}}} = \frac{\iint_{\mathcal{C}} dx dy}{\iint_{\mathcal{D}} dx dy} = \frac{\iint_{\mathcal{C}} p(x, y) dx dy}{\iint_{\mathcal{D}} p(x, y) dx dy} \\ &= \iint_{\mathcal{D}} \mathbb{1}\{(x, y) \in \mathcal{C}\} p(x, y) dx dy \\ &\approx \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{(x^{(i)}, y^{(i)}) \in \mathcal{C}\} = \frac{N_{\text{hit}}}{N} \end{aligned}$$

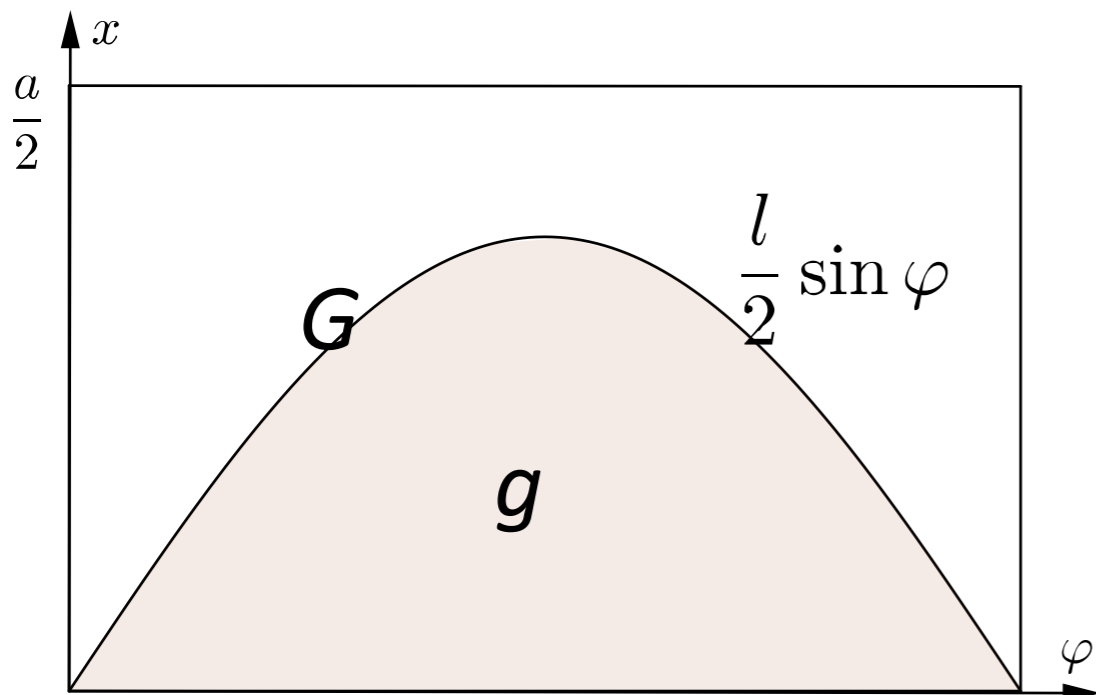
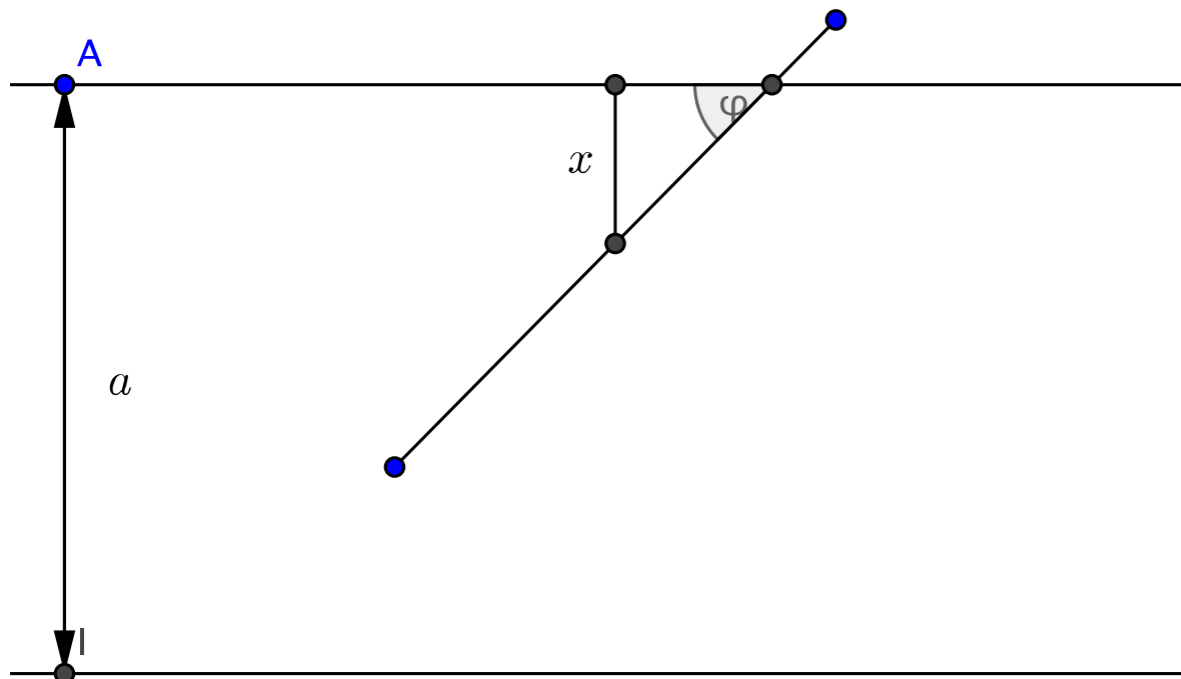


Another Way to Estimate π



$$\pi \approx \frac{2lN}{aN_{\text{hits}}}$$

Estimating π

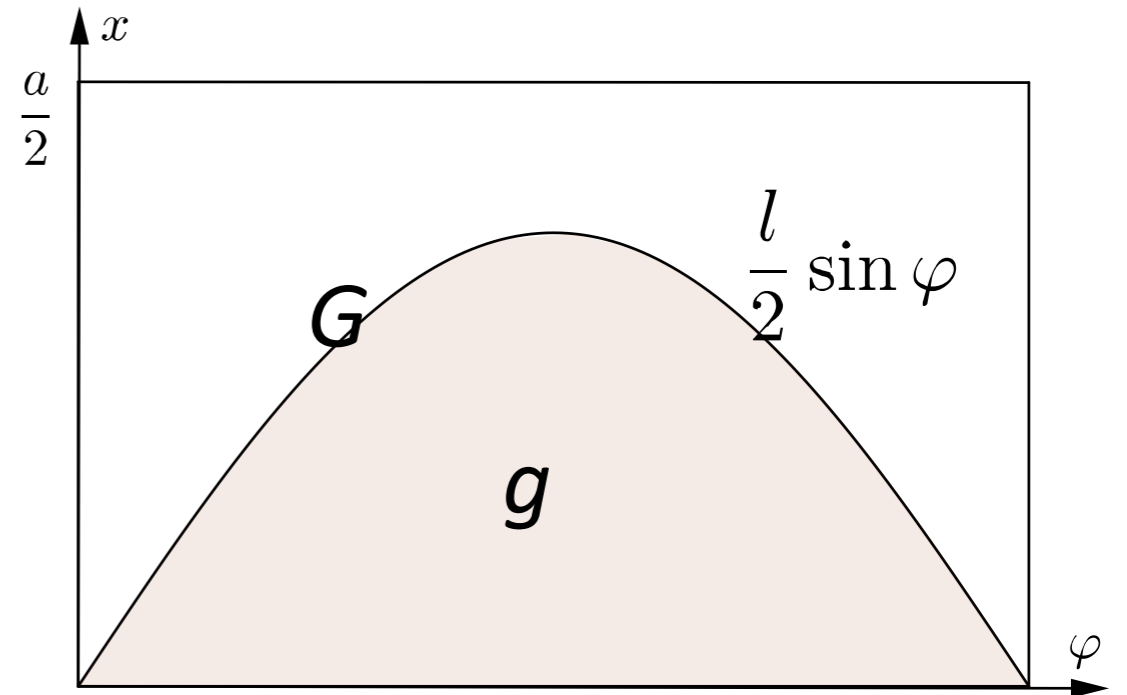


$$p_{\text{hit}} = \frac{S_g}{S_G} = \frac{\int_0^\pi \frac{l}{2} \sin \varphi \, d\varphi}{\frac{1}{2} a \pi} = \frac{2l}{\pi a}$$

$$p_{\text{hit}} \approx \frac{N_{\text{hits}}}{N} \quad \pi \approx \frac{2lN}{aN_{\text{hits}}}$$

Estimating π

$$\begin{aligned}\frac{2l}{\pi a} &= \frac{S_g}{S_G} = \frac{\iint_g dx dy}{\iint_G dx dy} \\ &= \frac{\iint_g p(x, y) dx dy}{\iint_G p(x, y) dx dy} \\ &= \iint_G \mathbb{1}\{(x, y) \in g\} p(x, y) dx dy \\ &= \mathbb{E}[\mathbb{1}\{(x, y) \in g\}] \\ &\approx \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{(x, y) \in g\} = \frac{N_{\text{hits}}}{N}\end{aligned}$$



Monte Carlo: Motivation

- Estimation
- Optimization
- Model Selection

Random Variable Generation

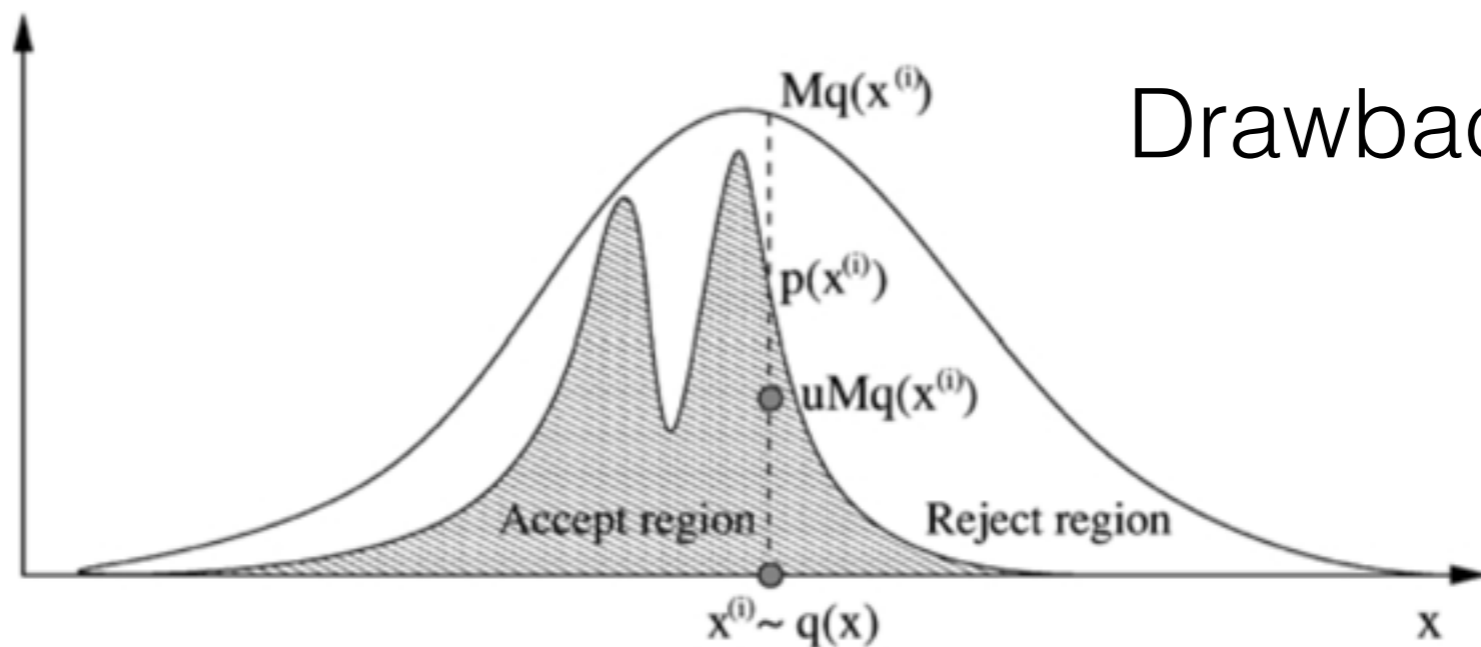
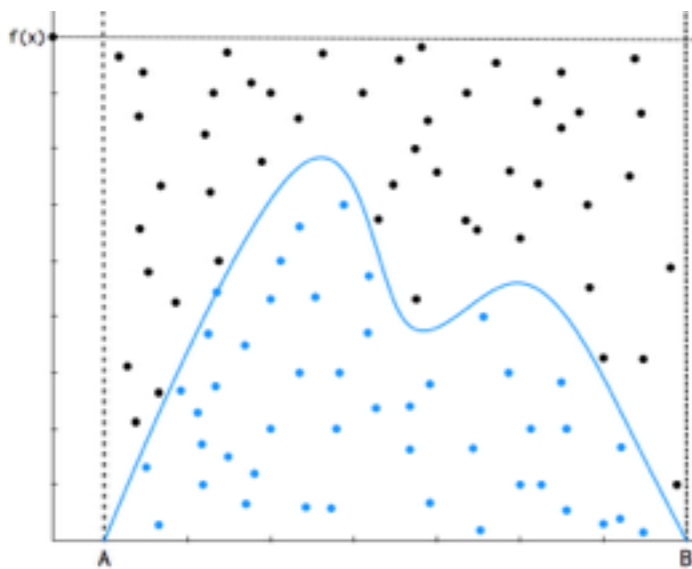
- if $U \sim \mathcal{U}_{[0,1]}$, then the random variable $F^{-1}(U)$ has the distribution F
where $F^{-1}(u) = \inf\{x : F(x) \geq u\}$
- The theorem above is theoretical beautiful, but unpractical in real application
- It is hard to compute the generalized inverse function
- High dimensional cases?

Rejection Sampling

Set $i = 1$

Repeat until $i = N$

1. Sample $x^{(i)} \sim q(x)$ and $u \sim \mathcal{U}(0,1)$.
2. If $u < \frac{p(x^{(i)})}{Mq(x^{(i)})}$ then accept $x^{(i)}$ and increment the counter i by 1. Otherwise, reject.

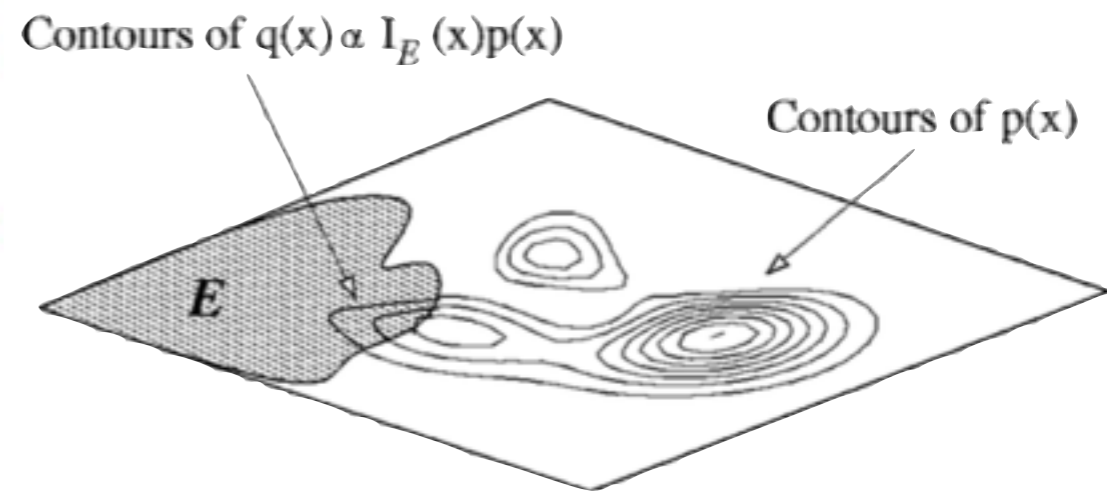
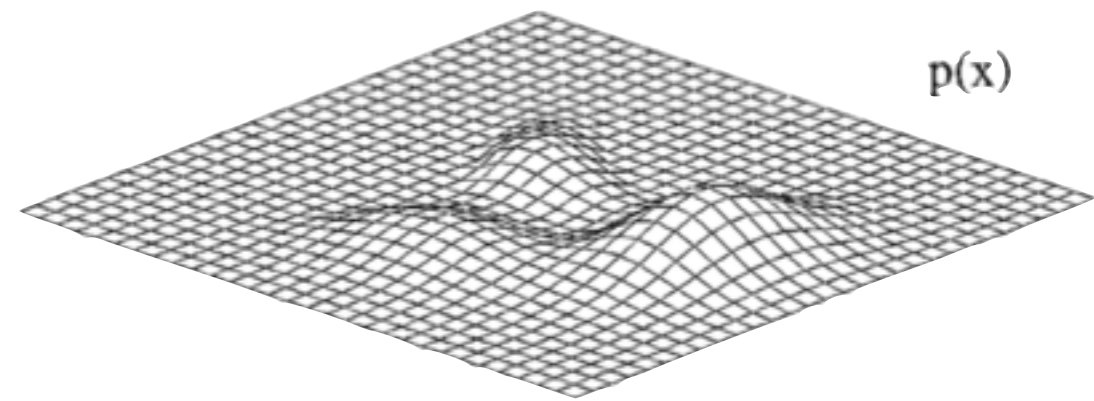
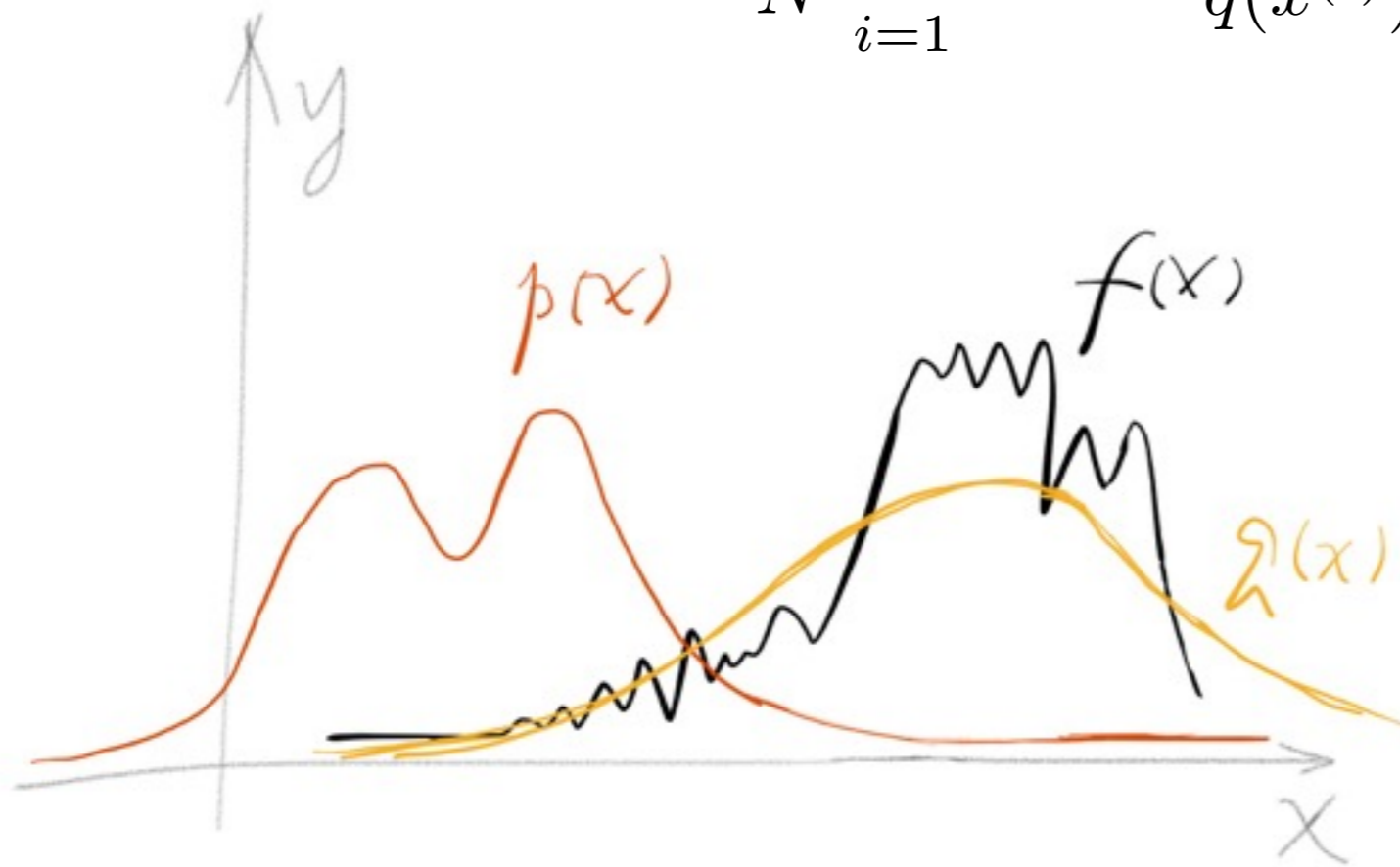


Drawbacks?

Importance Sampling

$$\int f(x)p(x) dx = \int f(x) \frac{p(x)}{q(x)} q(x) dx$$

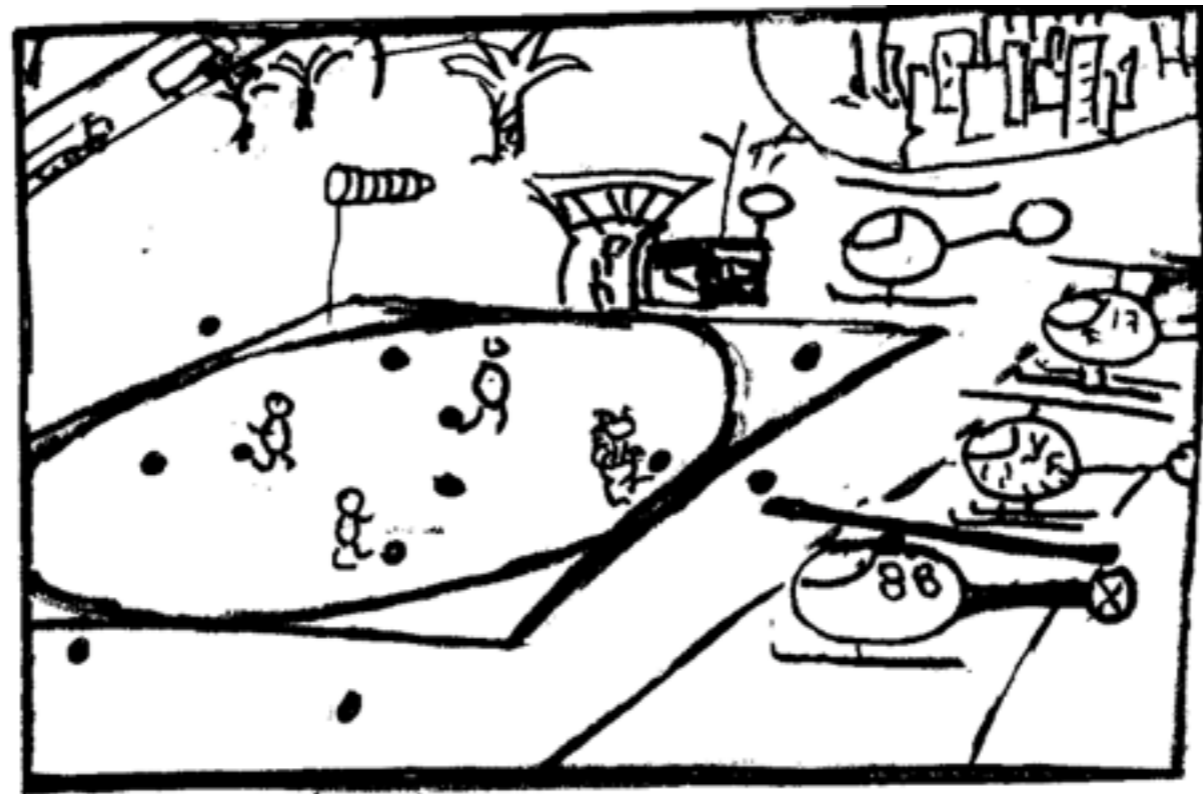
$$\approx \frac{1}{N} \sum_{i=1}^N f(x^{(i)}) \frac{p(x^{(i)})}{q(x^{(i)})}$$

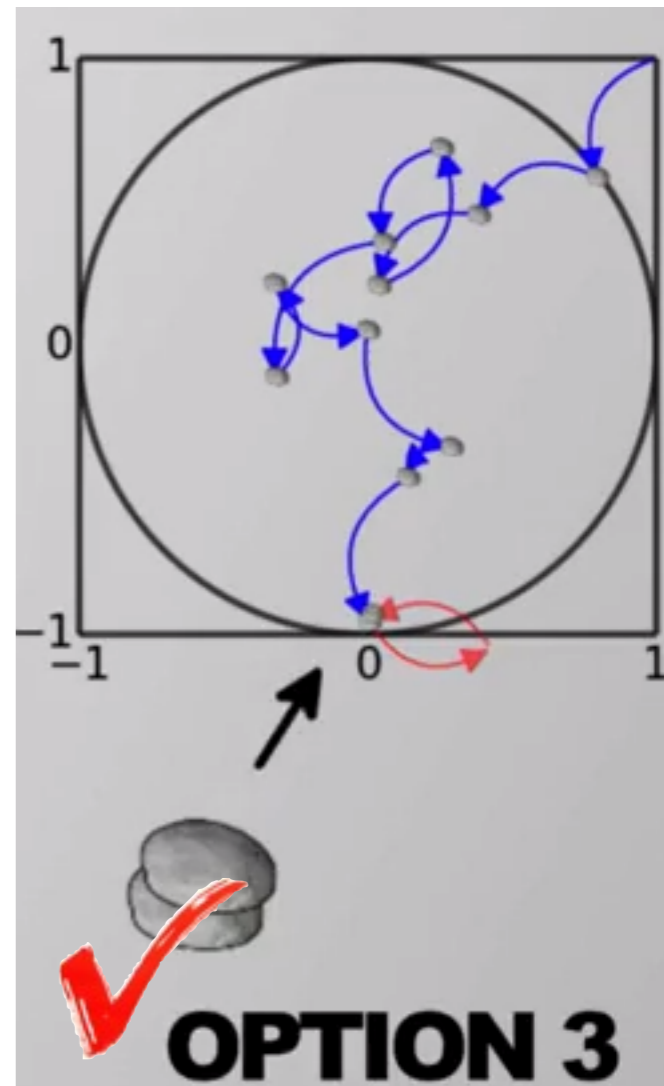
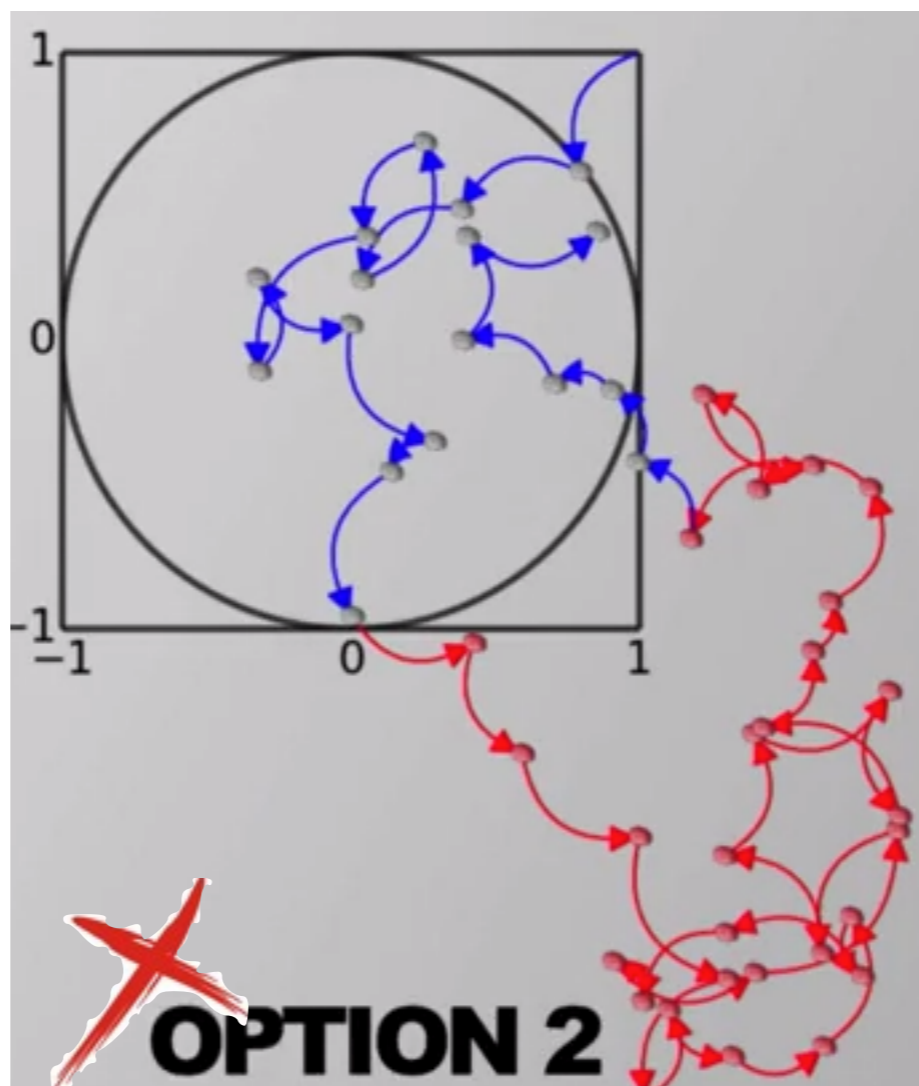
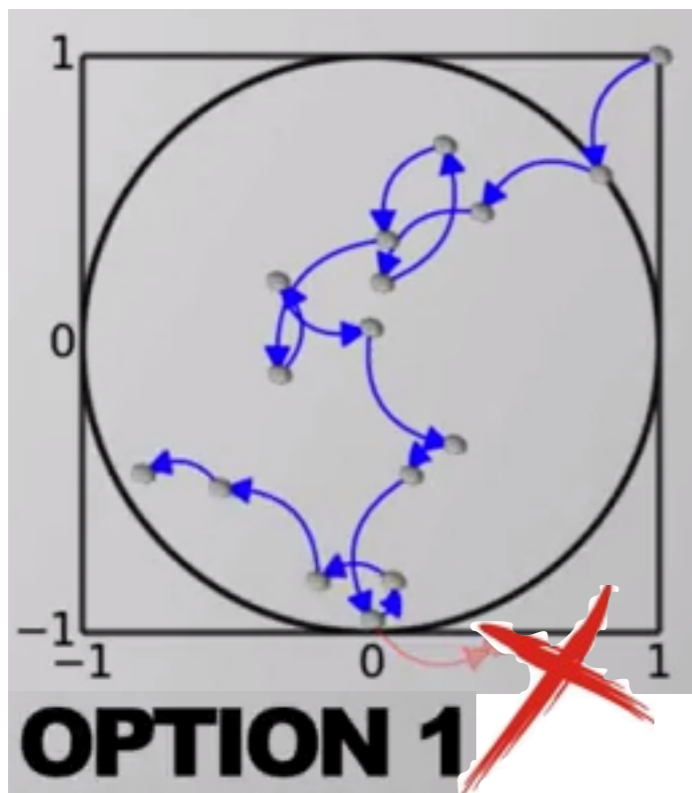
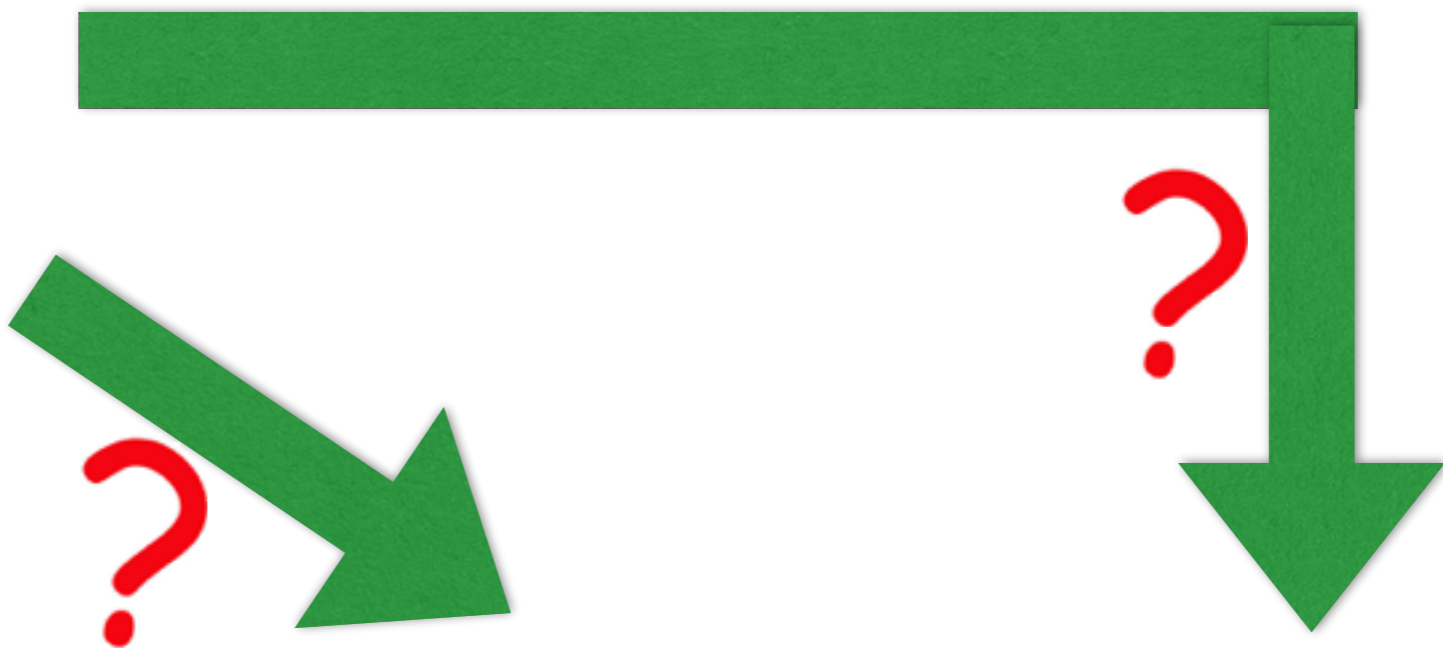
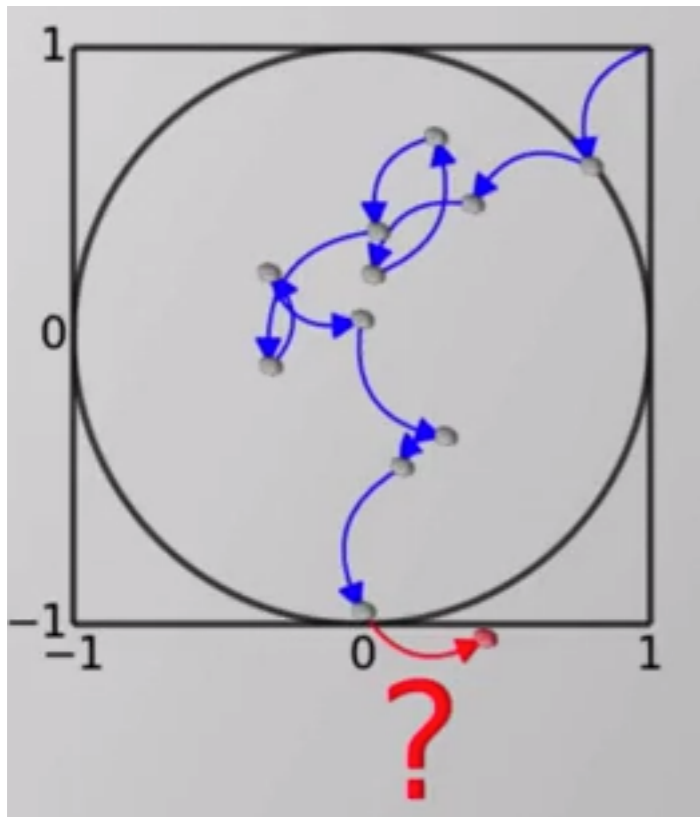


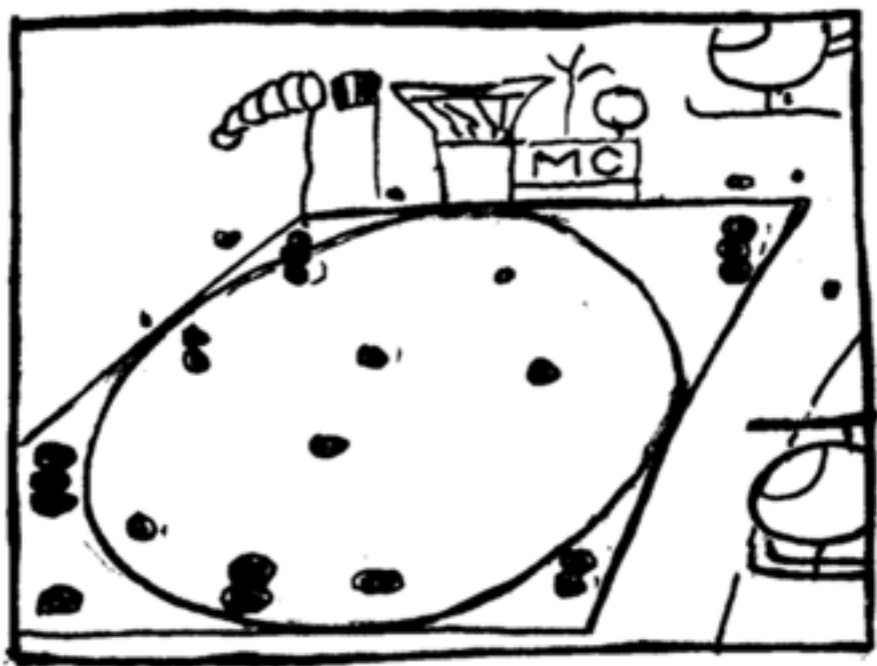
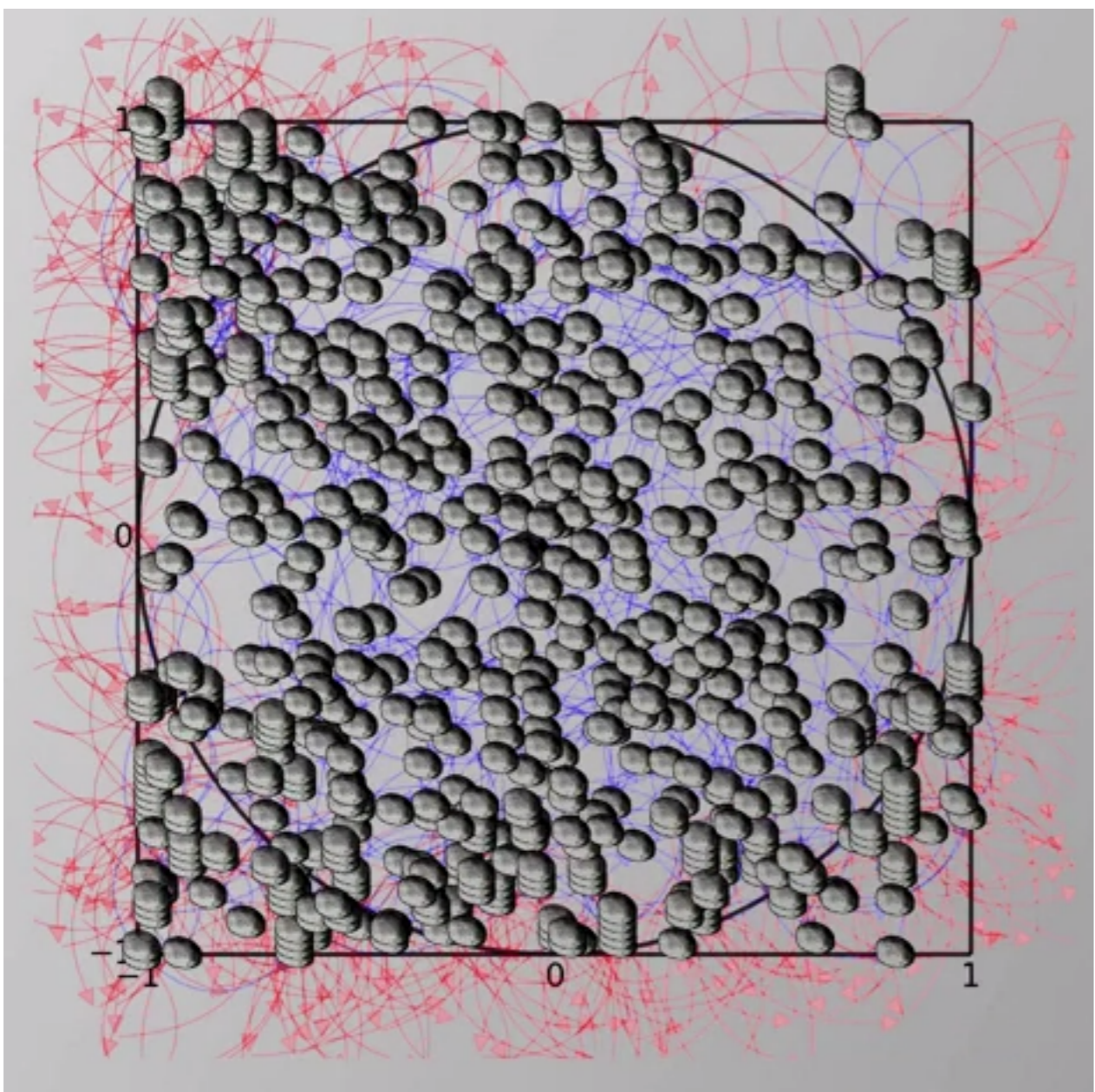
$$\text{Var}_p[f(x)] > \text{Var}_q\left[\frac{p(x)}{q(x)} f(x)\right]$$

- Does the i.i.d. property necessary for estimating an integral?

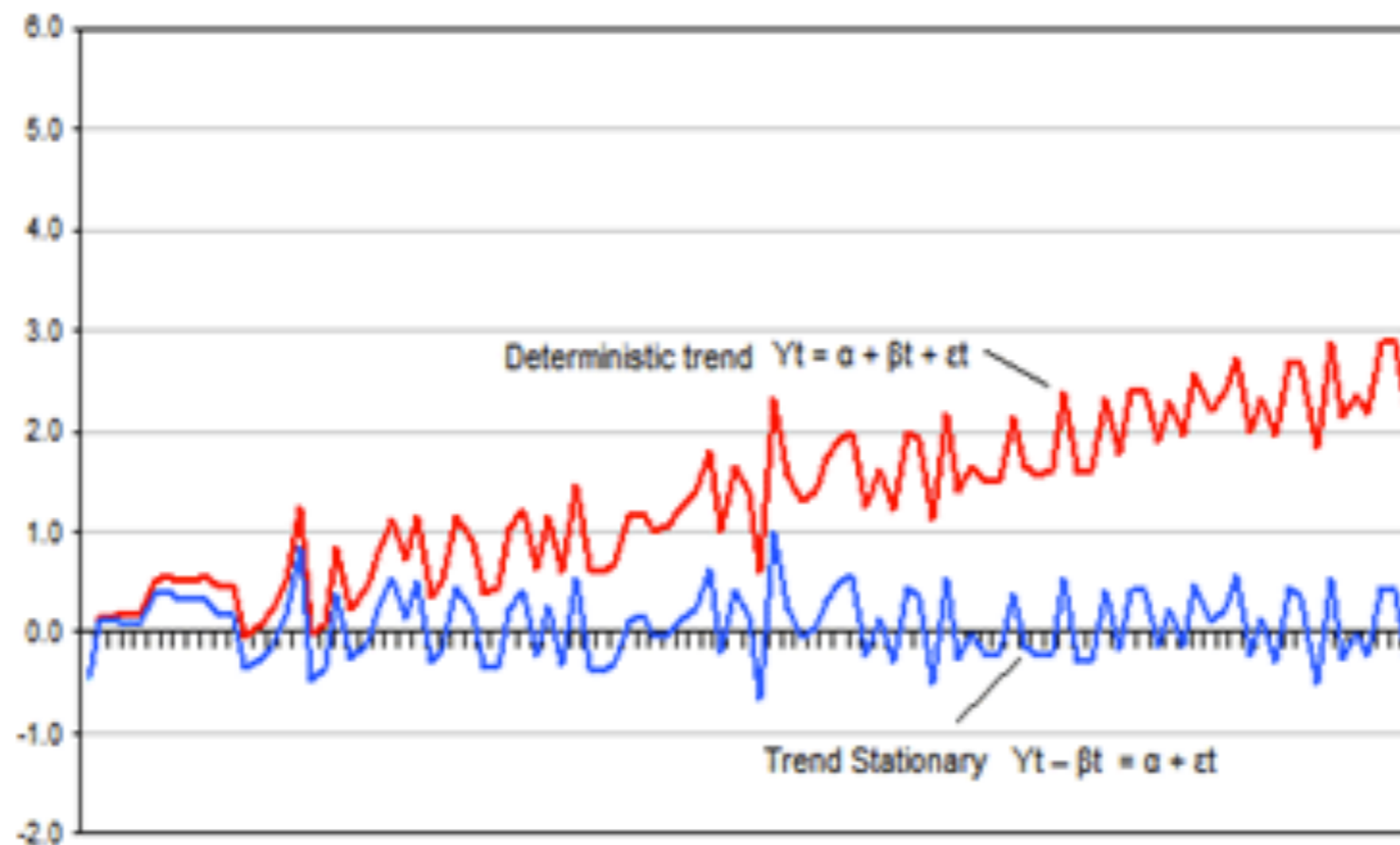
Another π Estimation Story







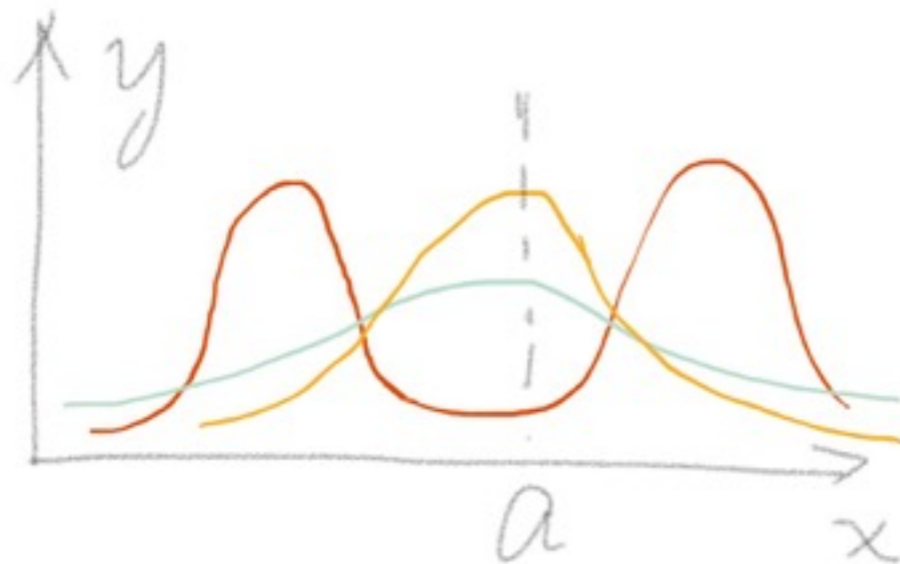
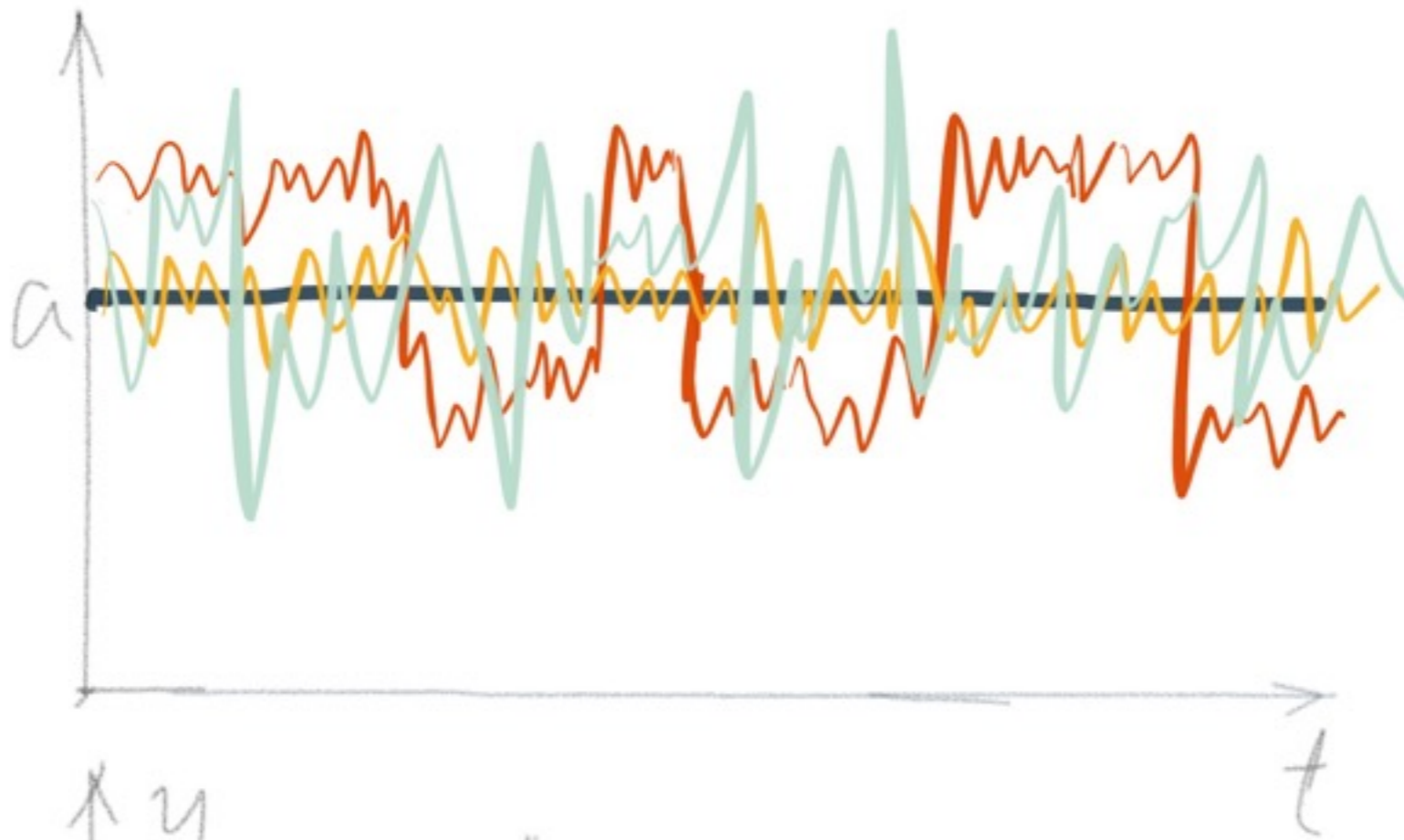
Stationary Stochastic Process



$$\mathbb{E}[X_t] = \text{constant} = \lim_{T \rightarrow +\infty} \frac{1}{T} \int_0^T X(t) dt$$

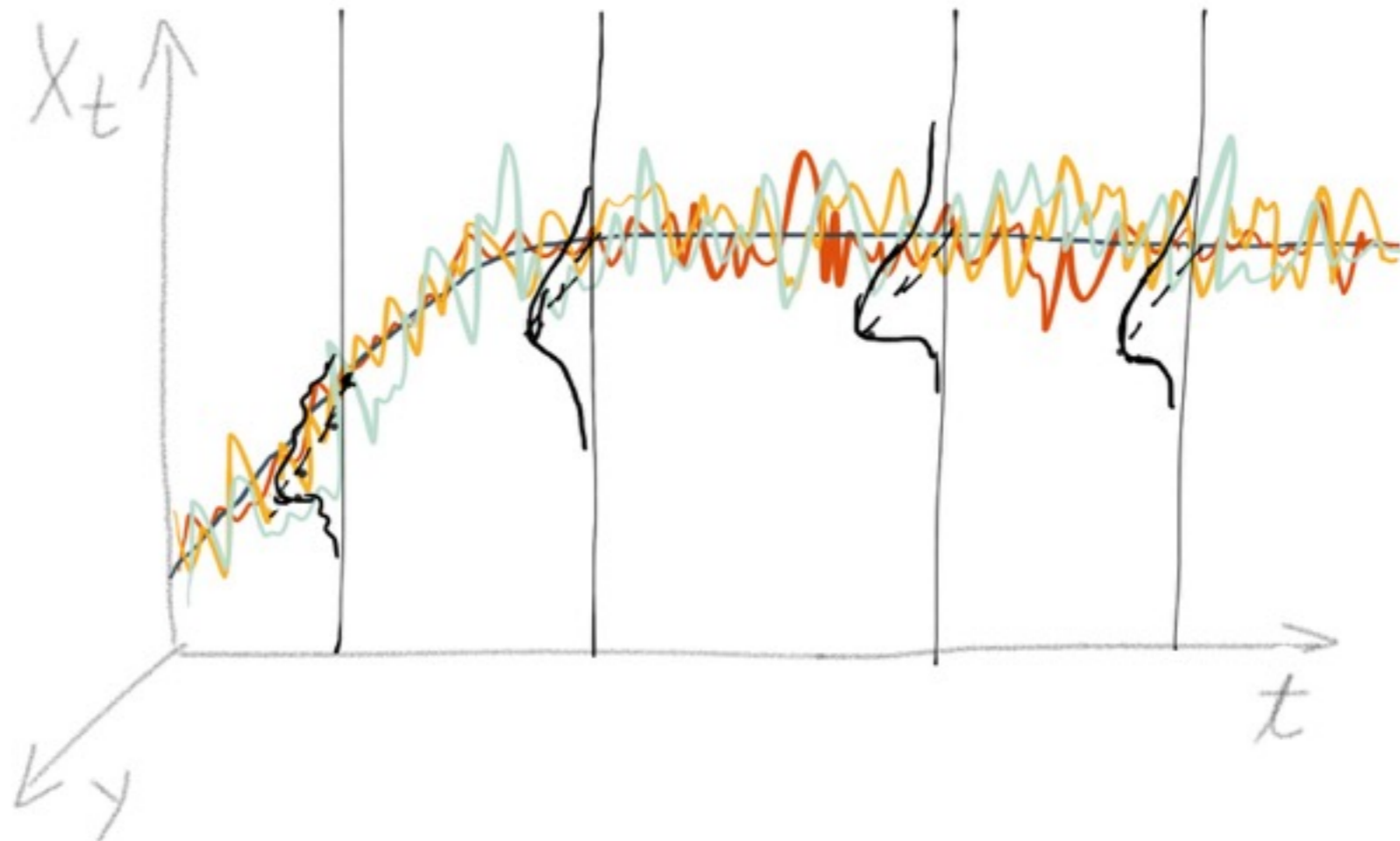
? ~~$$\mathbb{E}[f(X_t)] = \lim_{T \rightarrow +\infty} \frac{1}{T} \int_0^T f(X(t)) dt$$~~

Why $\mathbb{E}[f(X_t)] = \lim_{T \rightarrow +\infty} \frac{1}{T} \int_0^T f(X(t)) dt$ ~~?~~ ?

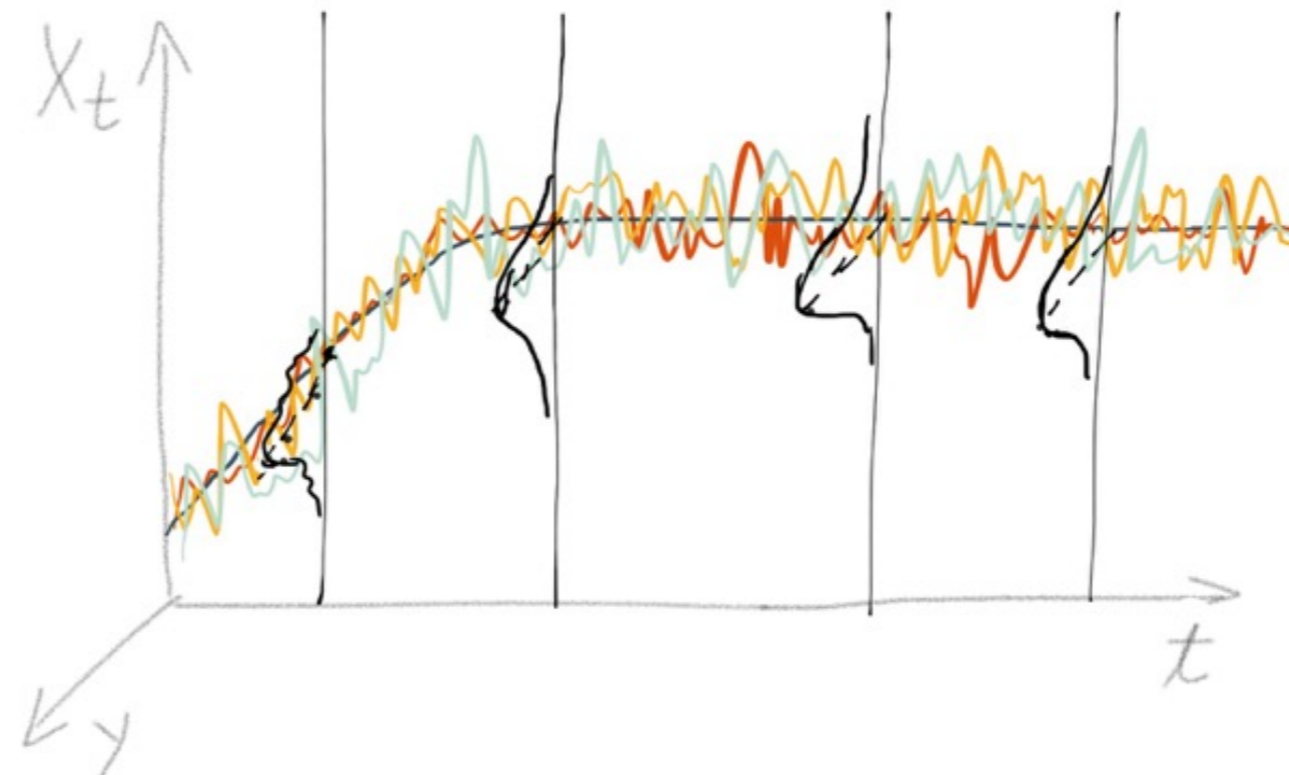
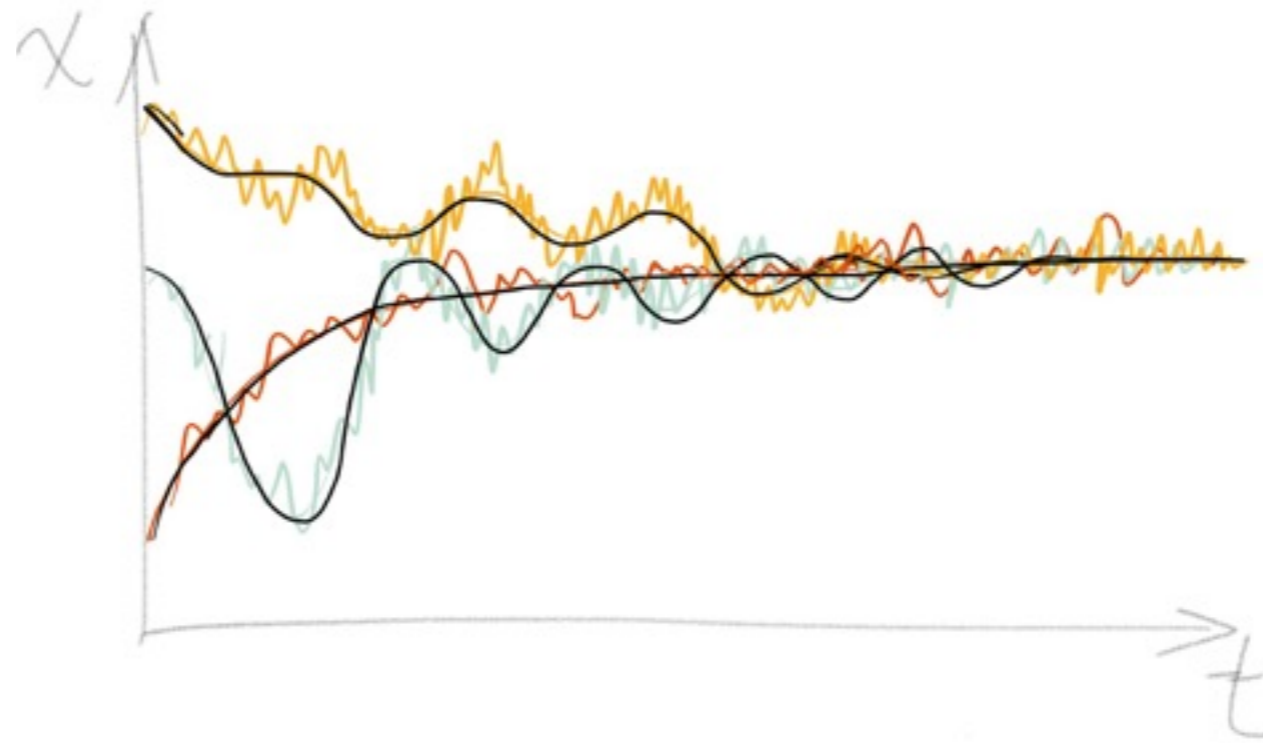


- How to make the equation below hold?

$$\mathbb{E}[f(X_t)] = \lim_{T \rightarrow +\infty} \frac{1}{T} \int_0^T f(X(t)) dt$$



Markov Chain Monte Carlo (MCMC): Intuitions



MCMC: Motivation

- Bayesian inference and learning

- Normalization
$$p(x|y) = \frac{p(y|x)p(x)}{\int_{\mathcal{X}} p(y|x')p(x') dx'}$$

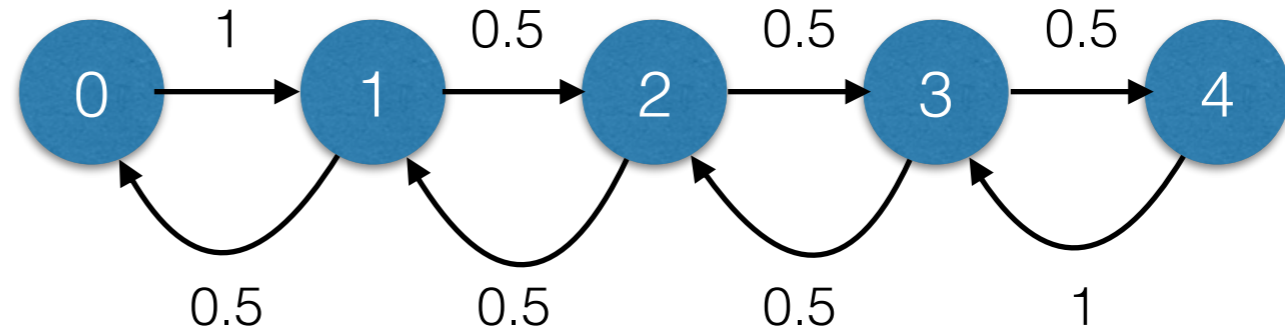
- Marginalization
$$p(x|y) = \int_{\mathcal{Z}} p(x, z, |y) dz$$

- Expectation
$$\mathbb{E}_{p(x|y)}[f(x)] = \int_{\mathcal{X}} f(x)p(x|y) dx$$

- Optimization
- Statistical mechanics: Compute the partition function
- Penalized likelihood model selection

MCMC: Important Theoretical Results

- Irreducibility
- Aperiodicity



- Detailed Balance Condition

$$\pi_i P_{ij} = \pi_j P_{ji}$$

MCMC: Applications to Sampling Algorithms: Metropolis-Hastings Alg. (MHA)

- Proposal Distribution

$$q(i, j) = q(x^{(j)} | x^{(i)})$$

- Acceptance Prob.

$$\alpha(i, j) = \mathcal{A}(x^{(i)}, x^{(j)})$$

- Target Distribution

$$\pi(i) = p(x^{(i)})$$

- Detailed Balance

$$\pi(i)q(i, j)\alpha(i, j) = \pi(j)q(j, i)\alpha(j, i) \leftarrow \alpha(i, j) = \min\left\{1, \frac{\pi(j)q(j, i)}{\pi(i)q(i, j)}\right\}$$

1. Initialise $x^{(0)}$.

2. For $i = 0$ to $N - 1$

– Sample $u \sim \mathcal{U}_{[0,1]}$.

– Sample $x^* \sim q(x^* | x^{(i)})$.

– If $u < \mathcal{A}(x^{(i)}, x^*) = \min\left\{1, \frac{p(x^*)q(x^{(i)} | x^*)}{p(x^{(i)})q(x^* | x^{(i)})}\right\}$

$$x^{(i+1)} = x^*$$

else

$$x^{(i+1)} = x^{(i)}$$

Special Cases of MHA

- Independent Sampler: Assuming that the proposal is independent of the current state

$$q(x^{(j)} | x^{(i)}) = q(x^{(j)})$$

$$\alpha(i, j) = \min \left\{ 1, \frac{p(x^{(j)})q(x^{(i)})}{p(x^{(i)})q(x^{(j)})} \right\} = \min \left\{ 1, \frac{w(x^{(j)})}{w(x^{(i)})} \right\}$$

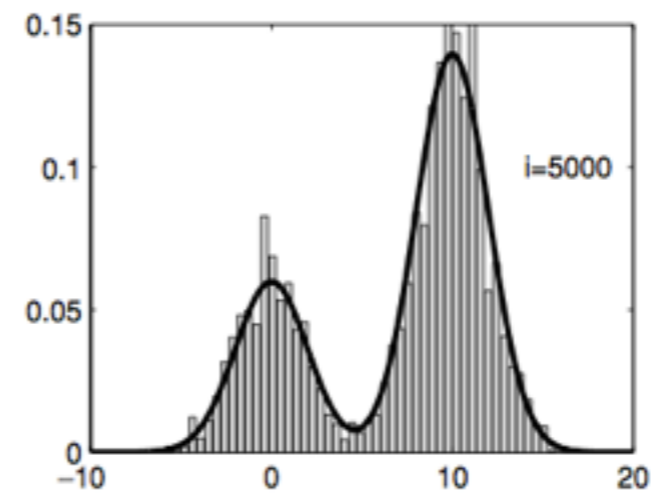
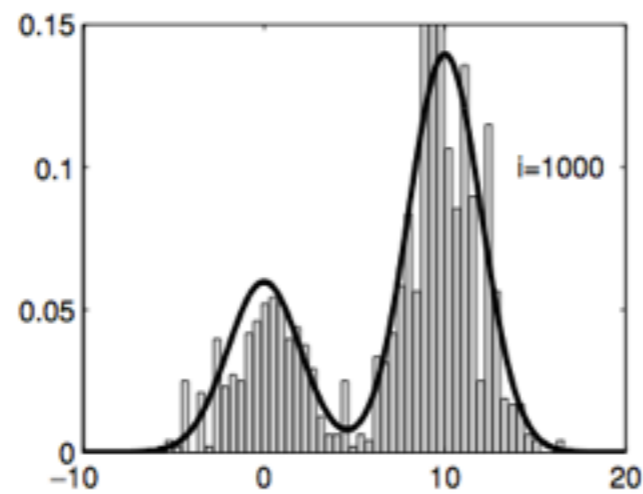
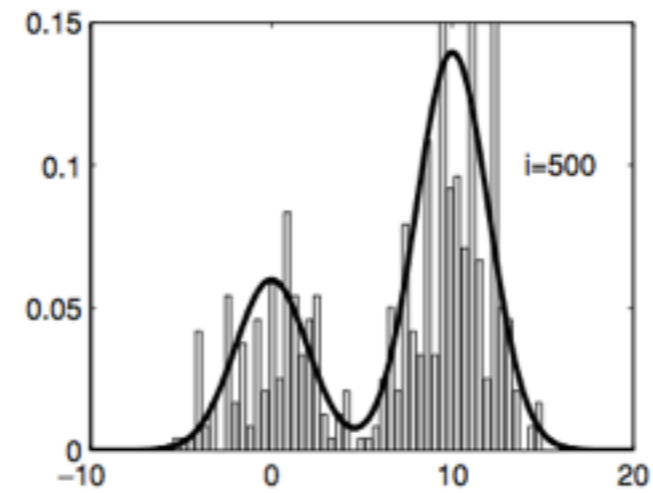
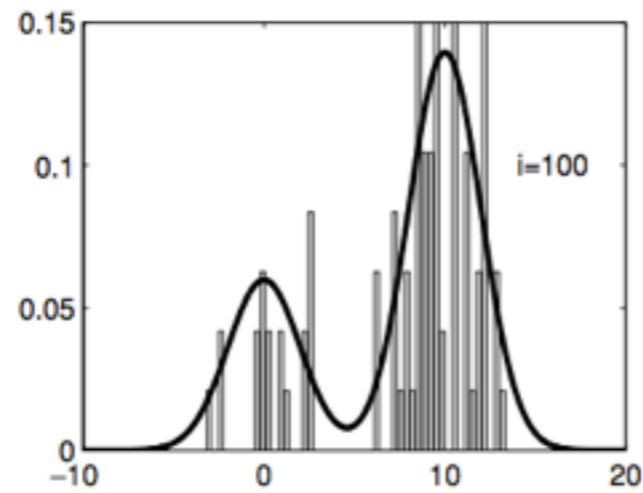
- Metropolis Sampler: Assuming that the proposal is symmetric

$$q(x^{(j)} | x^{(i)}) = q(x^{(i)} | x^{(j)})$$

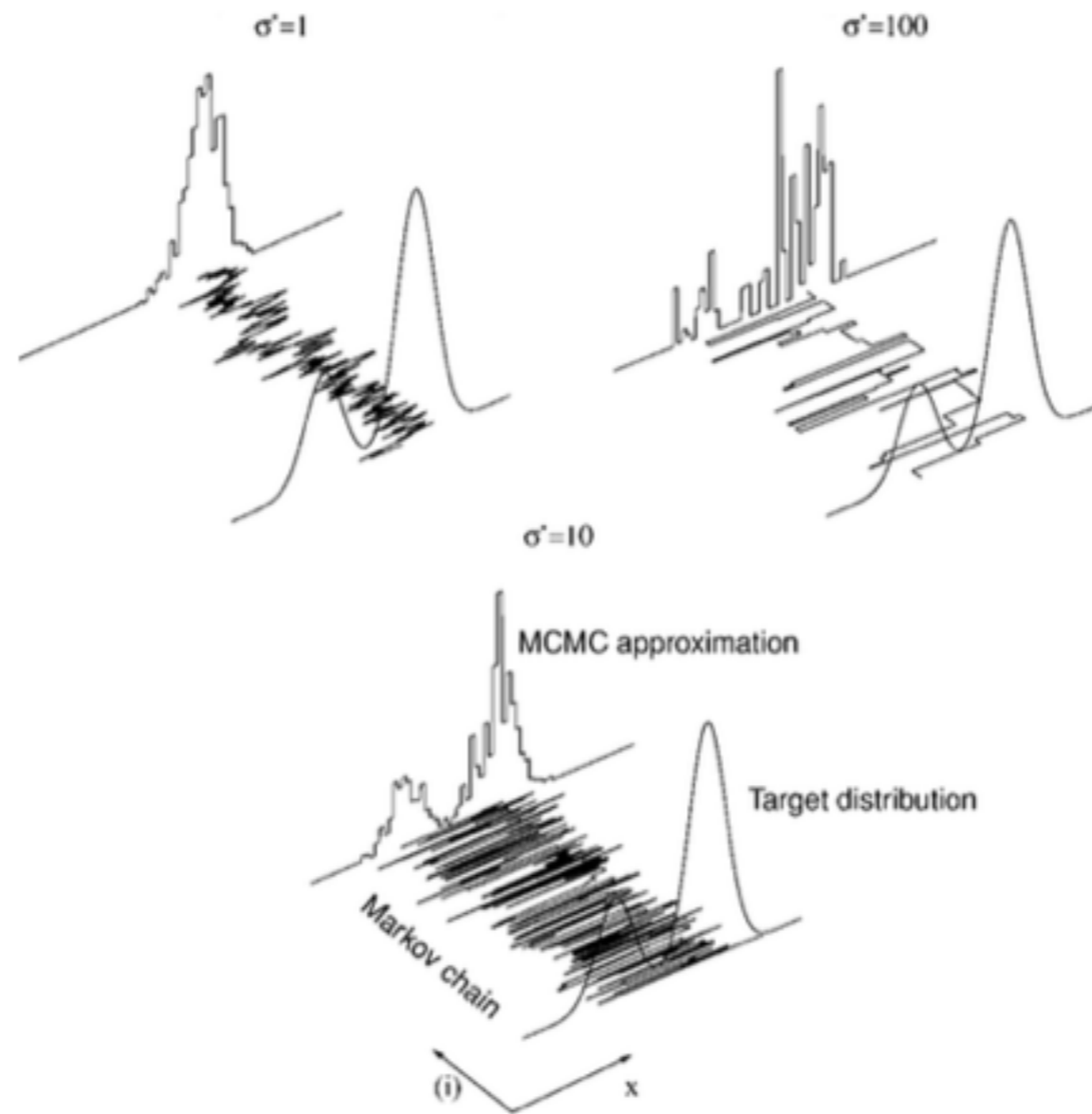
$$\alpha(i, j) = \min \left\{ 1, \frac{p(x^{(j)})}{p(x^{(i)})} \right\}$$

- Gibbs Sampler: Good property: no rejections

MHA Running Examples



MHA Running Examples



Simulated Annealing (SA)

- **Annealing** is the process of heating a solid until thermal stresses are released. Then, in cooling it very slowly to the ambient temperature until perfect crystals emerge. The quality of the results strongly depends on the cooling temperature. The final state can be interpreted as an energy state (crystalline potential energy) which is **lowest** if a perfectly crystal emerged

$$\Pr(\mathbf{x}) = \frac{1}{Z} \exp(-E(\mathbf{x})/k_B T)$$

$$\Pr(\mathbf{x} \rightarrow \mathbf{x}_p) = \begin{cases} 1, & \Delta E < 0; \\ \exp(-\Delta E/k_B T), & \Delta E \geq 0. \end{cases}$$

Simulated Annealing Algorithm

1. Initialize $\mathbf{x}^{(0)}$ and set $T_0, k = 0$.
2. For $i = 0$ to $N - 1$
 - Loop until converge
 - Sample $u \sim \mathcal{U}_{[0,1]}$.
 - Generate candidate state \mathbf{x}^* according to current state $\mathbf{x}^{(k)}$
 - $\Delta E = E(\mathbf{x}^*) - E(\mathbf{x}^{(k)})$
 - Accept state \mathbf{x}^* with probability:

$$\Pr(\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*) = \begin{cases} 1, & \Delta E < 0; \\ \exp(-\Delta E/k_B T_i), & \Delta E \geq 0. \end{cases}$$

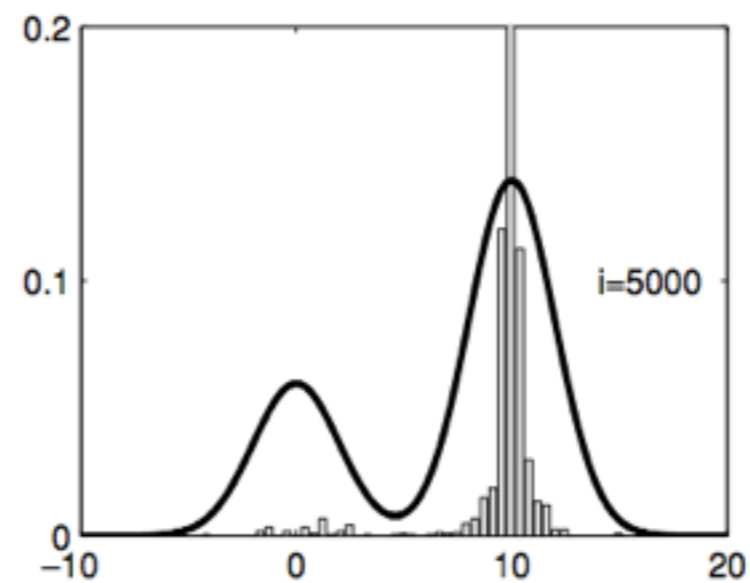
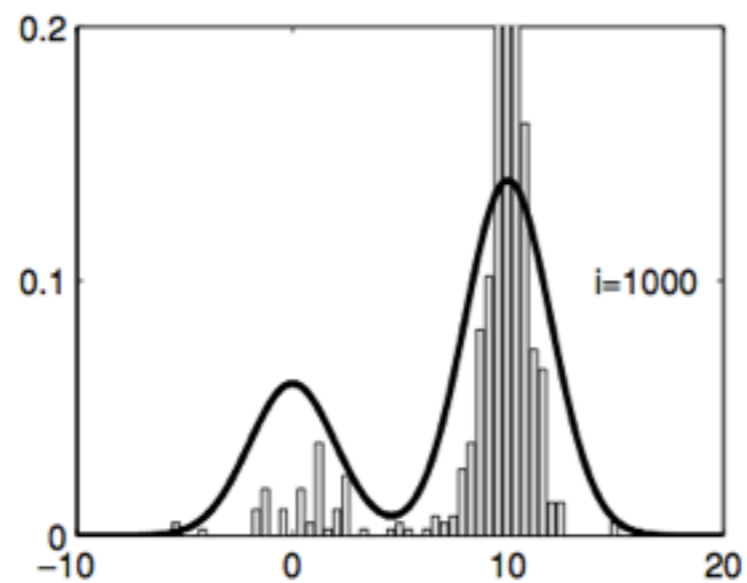
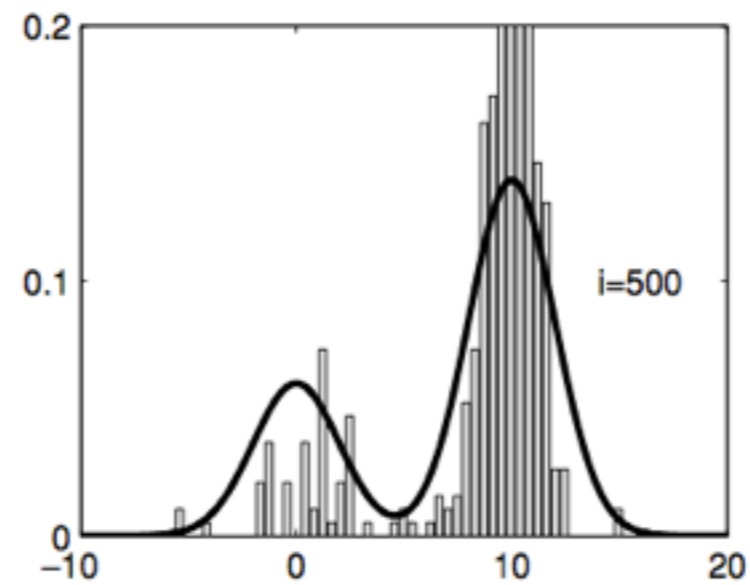
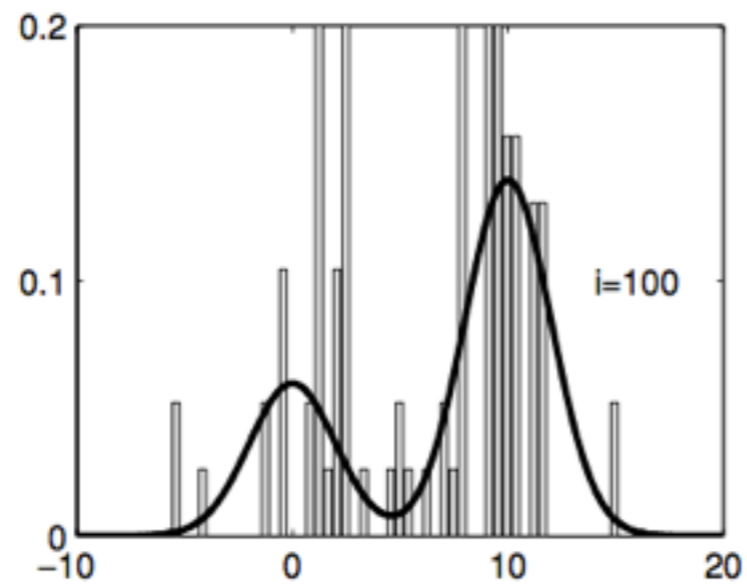
if $u < \Pr(\mathbf{x} \rightarrow \mathbf{x}^*)$, set $\mathbf{x}^{(k+1)} = \mathbf{x}^*$, else set $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$

- Set T_{i+1} according to cooling schedule.

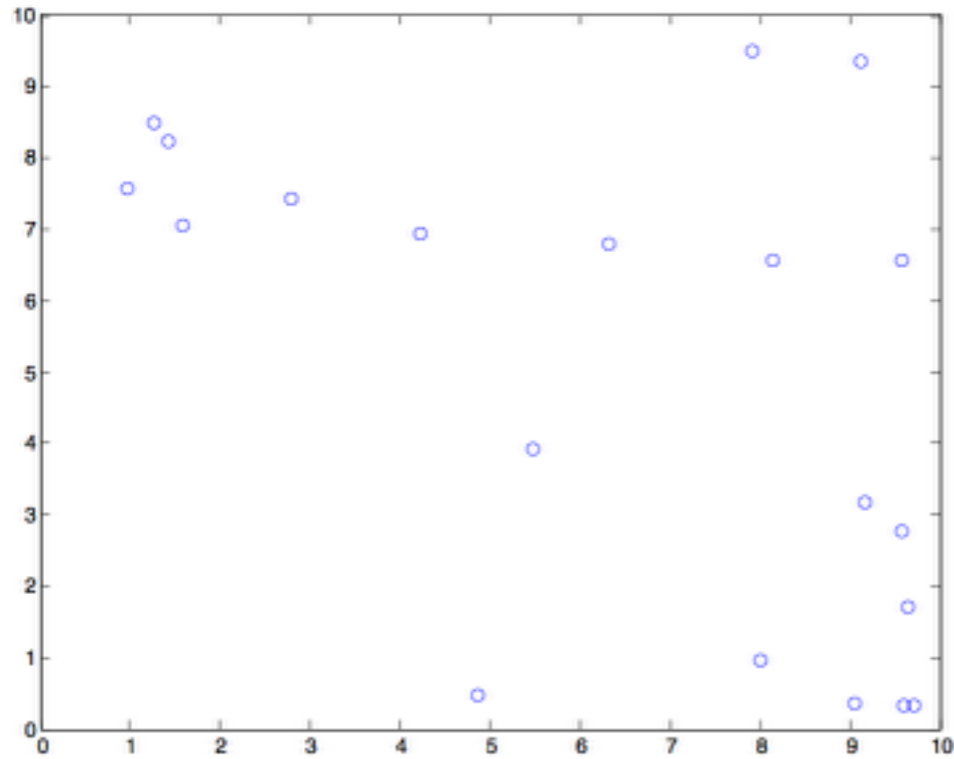
Simulated Annealing Algorithm

1. Initialise $x^{(0)}$ and set $T_0 = 1$.
2. For $i = 0$ to $N - 1$
 - Sample $u \sim \mathcal{U}_{[0,1]}$.
 - Sample $x^* \sim q(x^*|x^{(i)})$.
 - If $u < \mathcal{A}(x^{(i)}, x^*) = \min \left\{ 1, \frac{p^{\frac{1}{T_i}}(x^*)q(x^{(i)}|x^*)}{p^{\frac{1}{T_i}}(x^{(i)})q(x^*|x^{(i)})} \right\}$
 - $x^{(i+1)} = x^*$
 - else
 - $x^{(i+1)} = x^{(i)}$
 - Set T_{i+1} according to a chosen cooling schedule.

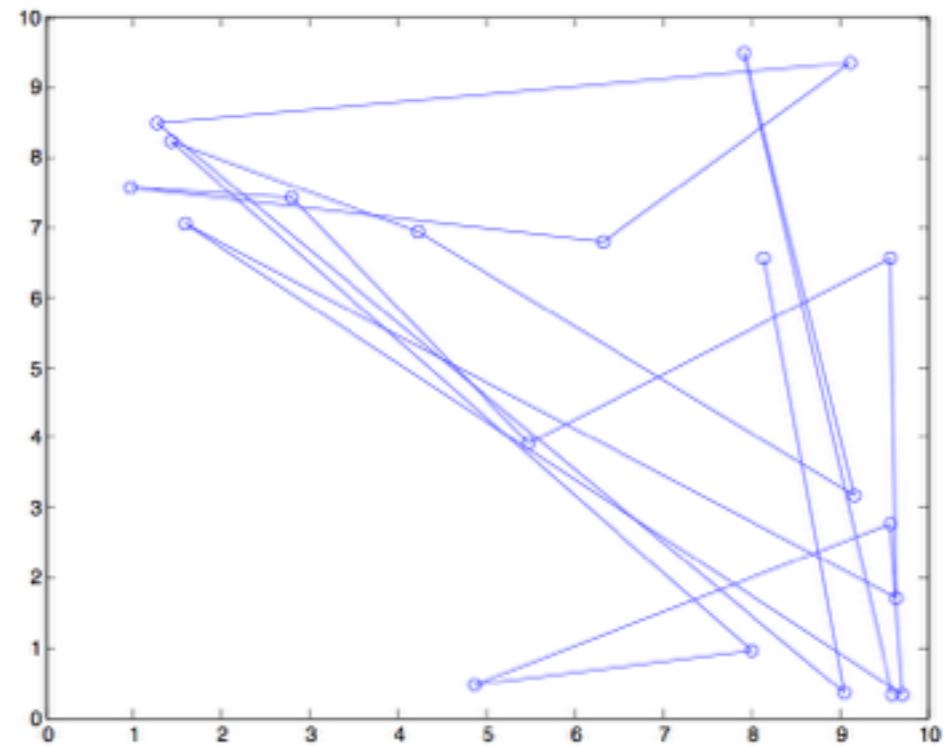
Toy Example: Maximization



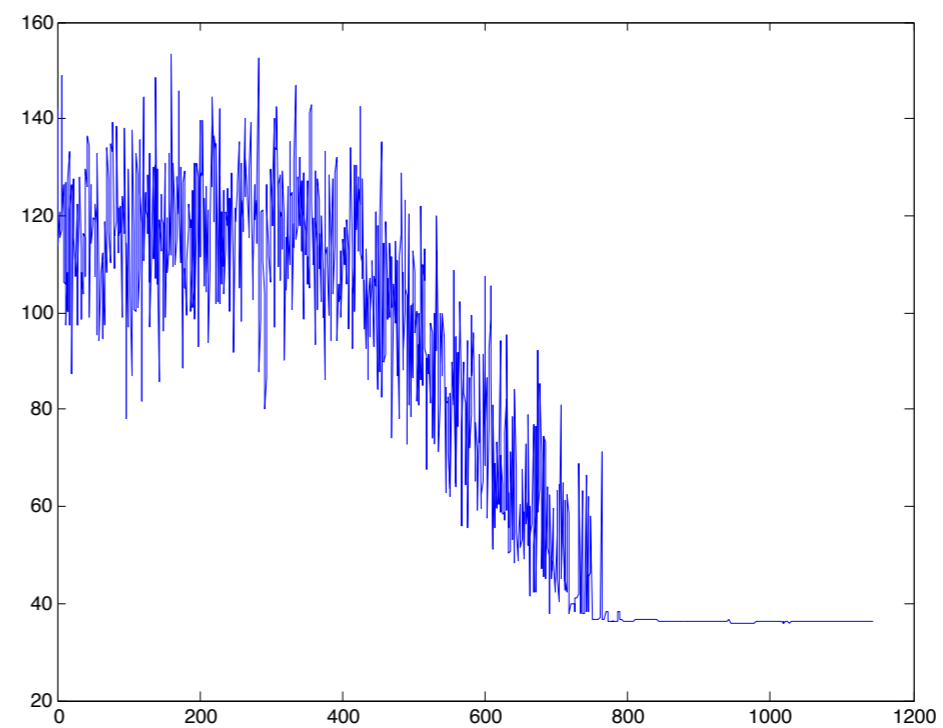
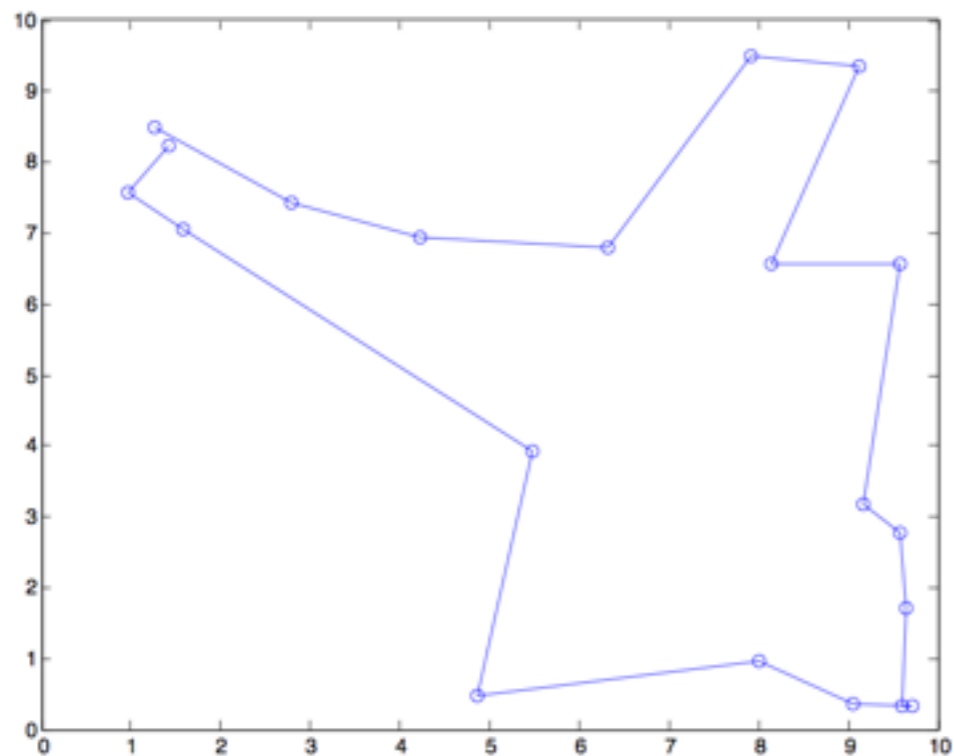
Solve Traveling Salesman Problem with SA



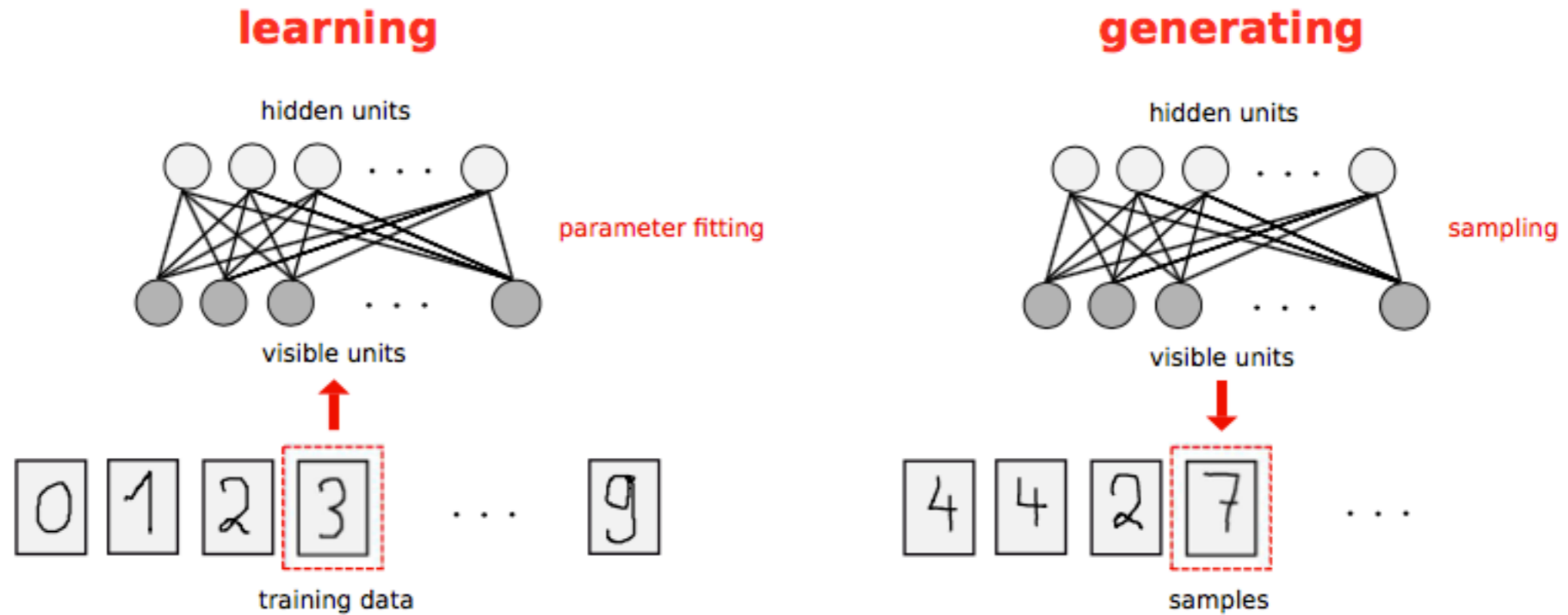
20 cities



Initial configuration

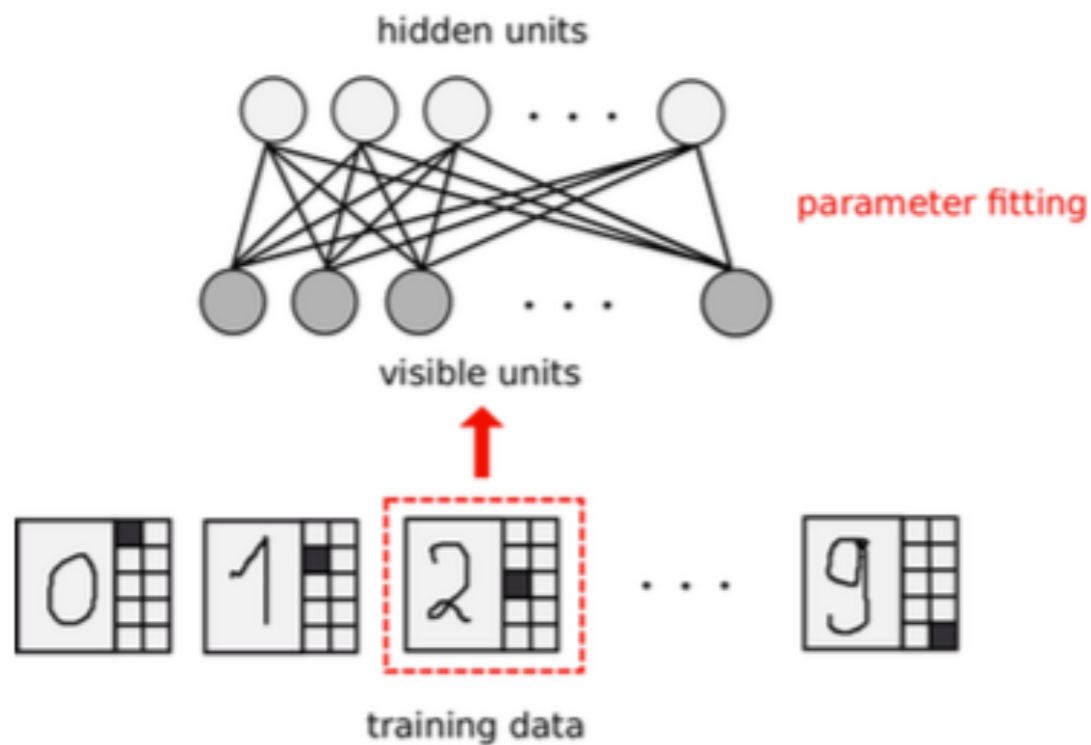


Restricted Boltzmann Machine (RBM)

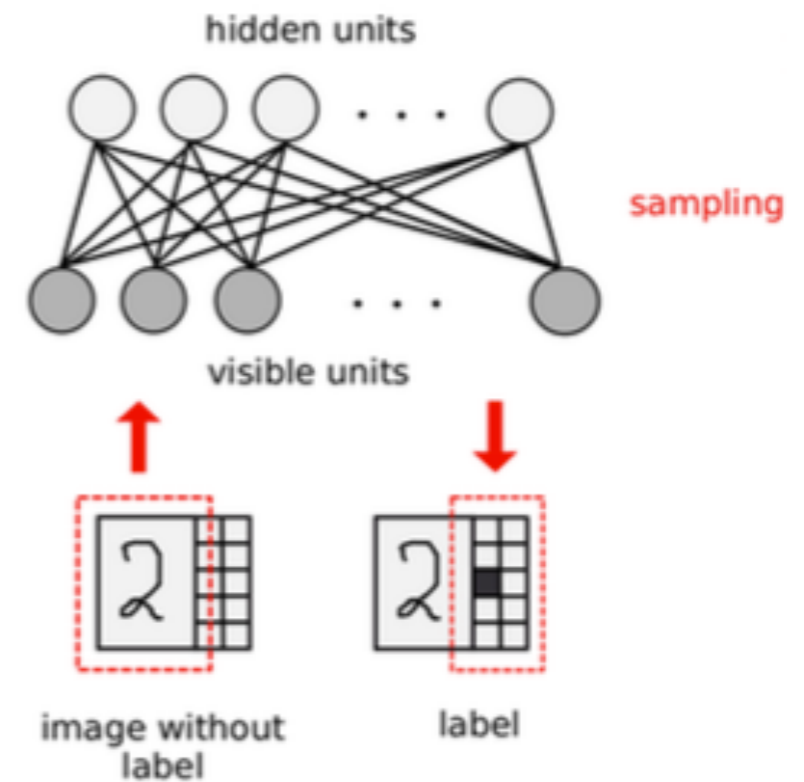


Restricted Boltzmann Machine (RBM)

learning with labels

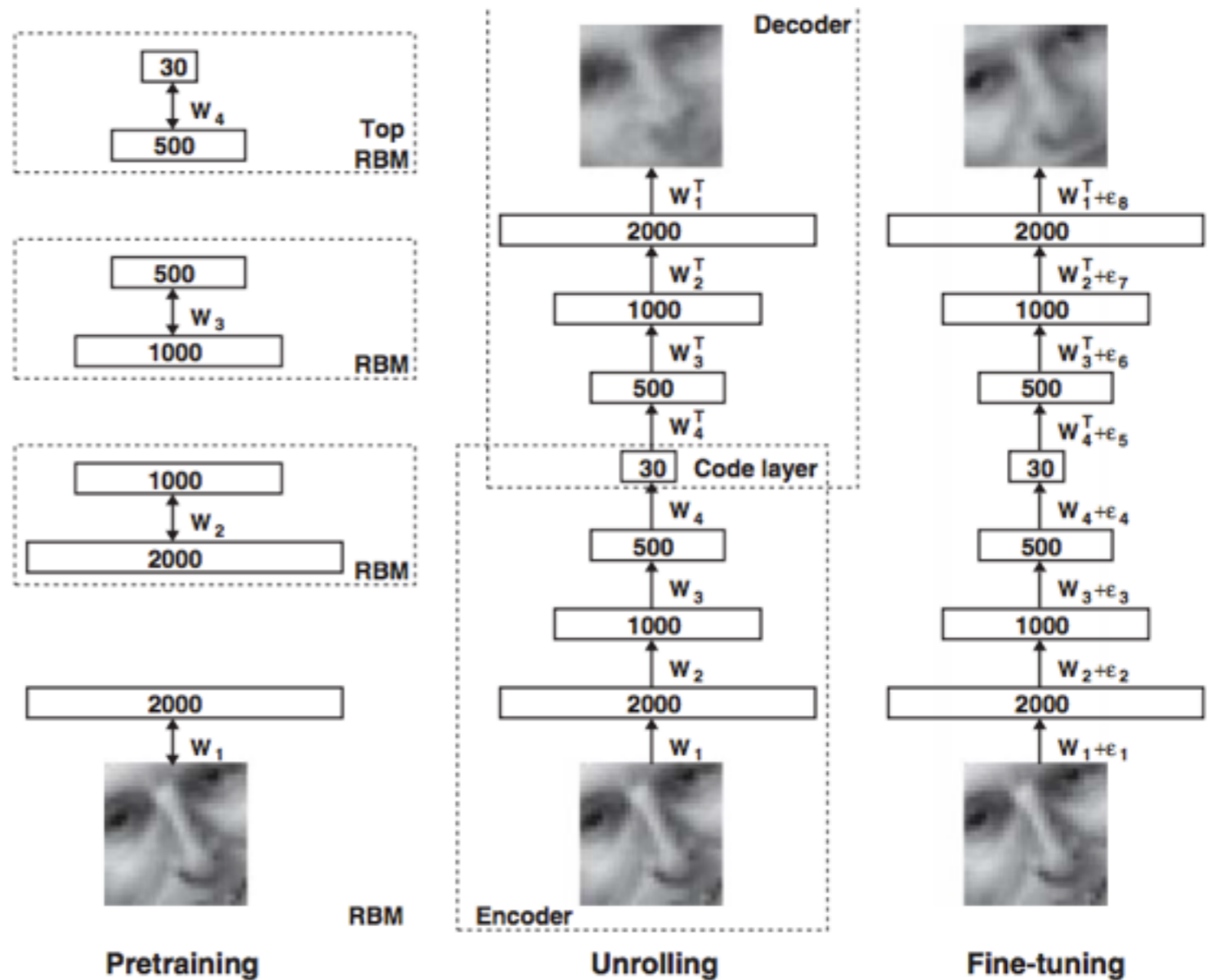
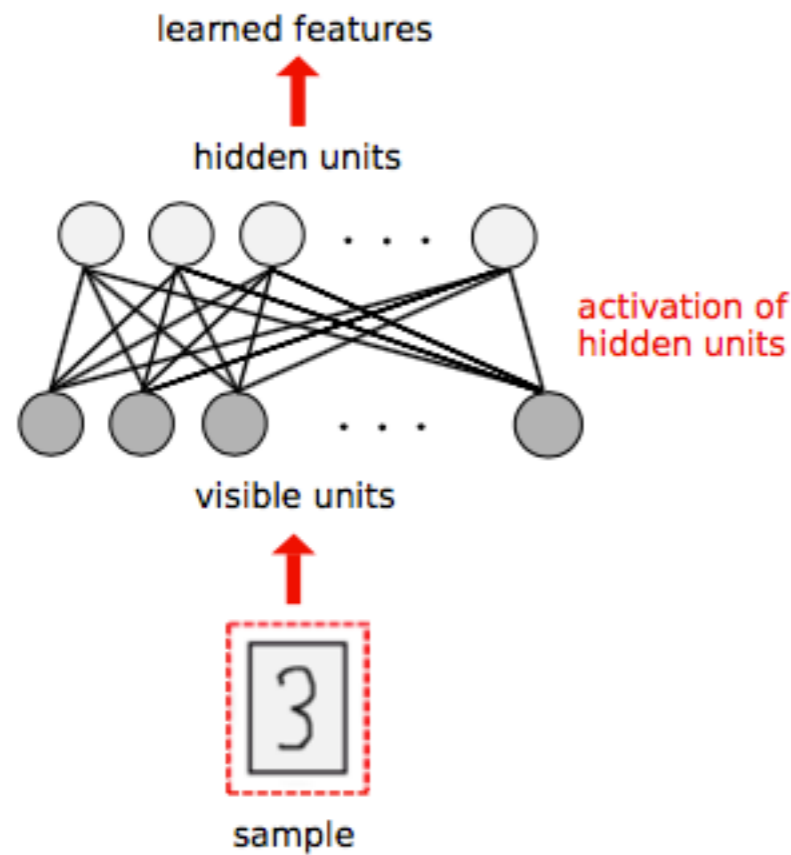


classification

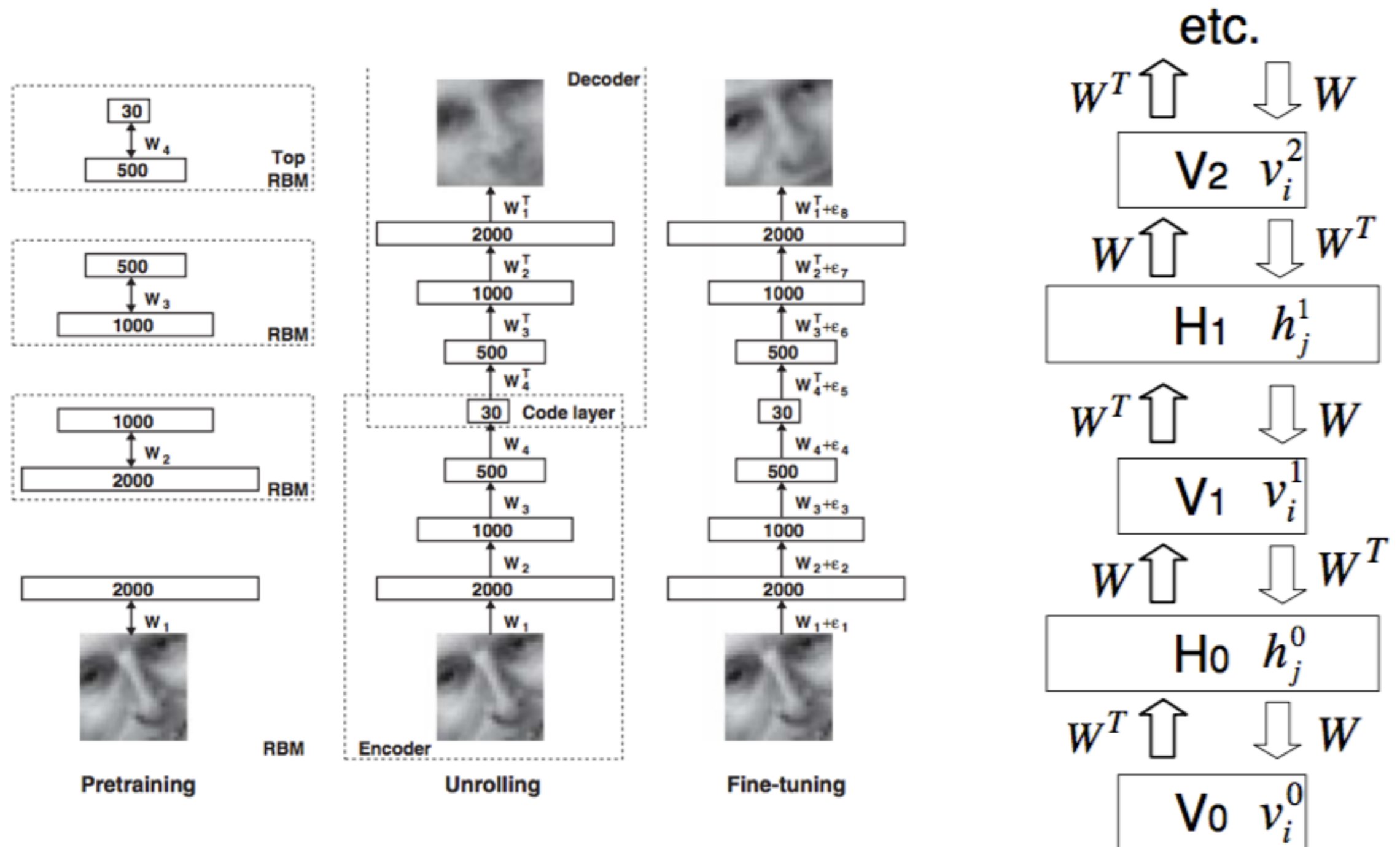


Restricted Boltzmann Machine (RBM)

feature mapping



Restricted Boltzmann Machine (RBM)



Basics to Learn RBM

- Fitting $p(\mathbf{v})$ to data

$$p(\mathbf{v}) = \frac{1}{Z} \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))$$

- Notations

$$\begin{aligned} E(\mathbf{v}, \mathbf{h}) &= -\mathbf{v}^T \mathbf{W} \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} \\ &= -\sum_{i=1}^m \sum_{j=1}^n w_{ij} h_i v_j - \sum_{j=1}^n b_j v_j - \sum_{i=1}^m c_i h_i \\ &= -\mathbf{W} \circ \mathbf{A} - \mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} \\ &= -\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}) \end{aligned}$$

$$\boldsymbol{\theta} = [\mathbf{W}(\cdot); \mathbf{b}; \mathbf{c}]$$

$$\boldsymbol{\phi}(\mathbf{x}) = [\mathbf{A}(\cdot); \mathbf{v}; \mathbf{h}]$$

$$\mathbf{W} = \begin{pmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1} & w_{m2} & \dots & w_{mn} \end{pmatrix}$$

$$\mathbf{A} = \begin{pmatrix} h_1 v_1 & h_1 v_2 & \dots & h_1 v_n \\ h_2 v_1 & h_2 v_2 & \dots & h_2 v_n \\ \vdots & \vdots & \ddots & \vdots \\ h_m v_1 & h_m v_2 & \dots & h_m v_n \end{pmatrix}$$

Learning RBM: Maximum Likelihood Estimation

- Given Observed Samples

$$S = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_\ell\}$$

- We can get

$$\frac{1}{\ell} \sum_{\mathbf{v} \in S} \frac{\partial \log p(\mathbf{v}; \boldsymbol{\theta})}{\partial w_{ij}} = \mathbb{E}_{p(\mathbf{h}|\mathbf{v}; \boldsymbol{\theta})q(\mathbf{v})} [h_i v_j] - \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} [h_i v_j]$$

Intractable

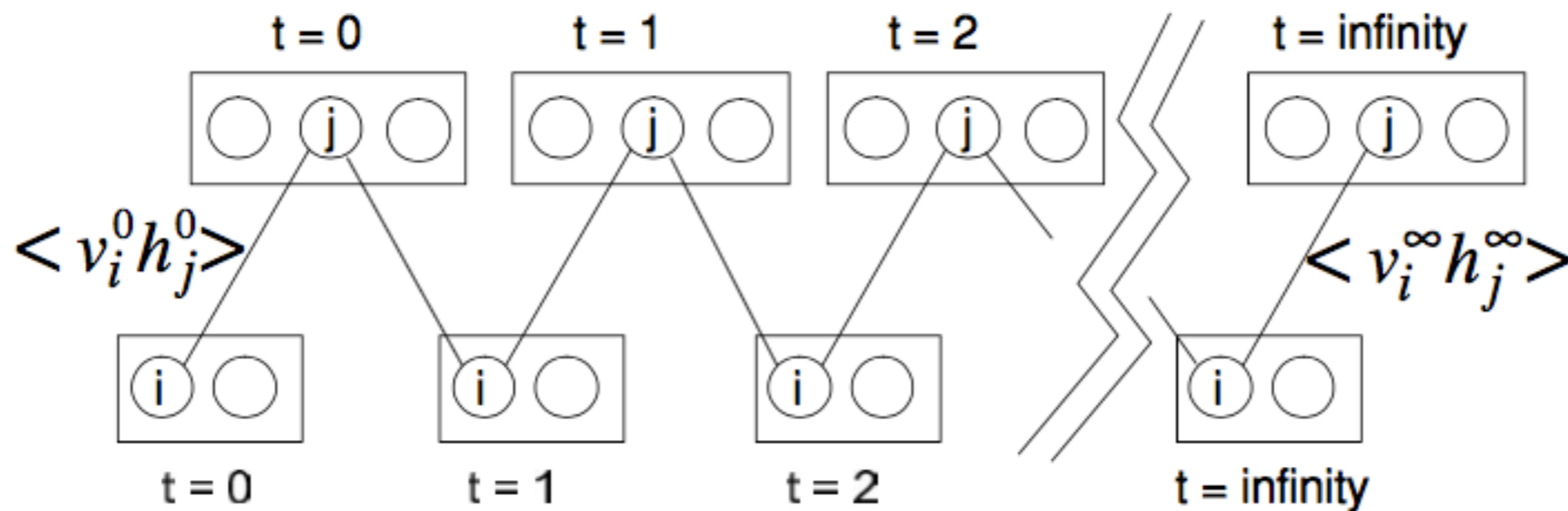
for

$$\frac{\partial \log p(\mathbf{v}_0; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_{p(\mathbf{h}|\mathbf{v}_0; \boldsymbol{\theta})} \boldsymbol{\phi}(\mathbf{v}_0, \mathbf{h}) - \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} [\boldsymbol{\phi}(\mathbf{x})]$$

$$\frac{\partial \log p(\mathbf{v}_0; \boldsymbol{\theta})}{\partial w_{ij}} = \mathbb{E}_{p(\mathbf{h}|\mathbf{v}_0; \boldsymbol{\theta})} [h_i v_j] - \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} [h_i v_j]$$

Contrastive Divergence

$$\frac{\partial \log p(\mathbf{v}_0; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_{p(\mathbf{h}|\mathbf{v}_0; \boldsymbol{\theta})} \boldsymbol{\phi}(\mathbf{v}_0, \mathbf{h}) - \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} [\boldsymbol{\phi}(\mathbf{x})]$$



$$\frac{\partial \log p(\mathbf{v}^{(0)}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \approx \mathbb{E}_{p(\mathbf{h}|\mathbf{v}^{(0)})} [\boldsymbol{\phi}(\mathbf{v}^{(0)}, \mathbf{h})] - \frac{1}{M} \sum_{i=1}^M \mathbb{E}_{p(\mathbf{h}|\mathbf{v}^{(k+i)})} [\boldsymbol{\phi}(\mathbf{v}^{(k+i)}, \mathbf{h})]$$

$$\text{CD}_k(\boldsymbol{\theta}, \mathbf{v}^{(0)}) = \mathbb{E}_{p(\mathbf{h}|\mathbf{v}^{(0)})} [\boldsymbol{\phi}(\mathbf{v}^{(0)}, \mathbf{h})] - \mathbb{E}_{p(\mathbf{h}|\mathbf{v}^{(k)})} [\boldsymbol{\phi}(\mathbf{v}^{(k)}, \mathbf{h})]$$

k-Step Contrastive Divergence Alg.

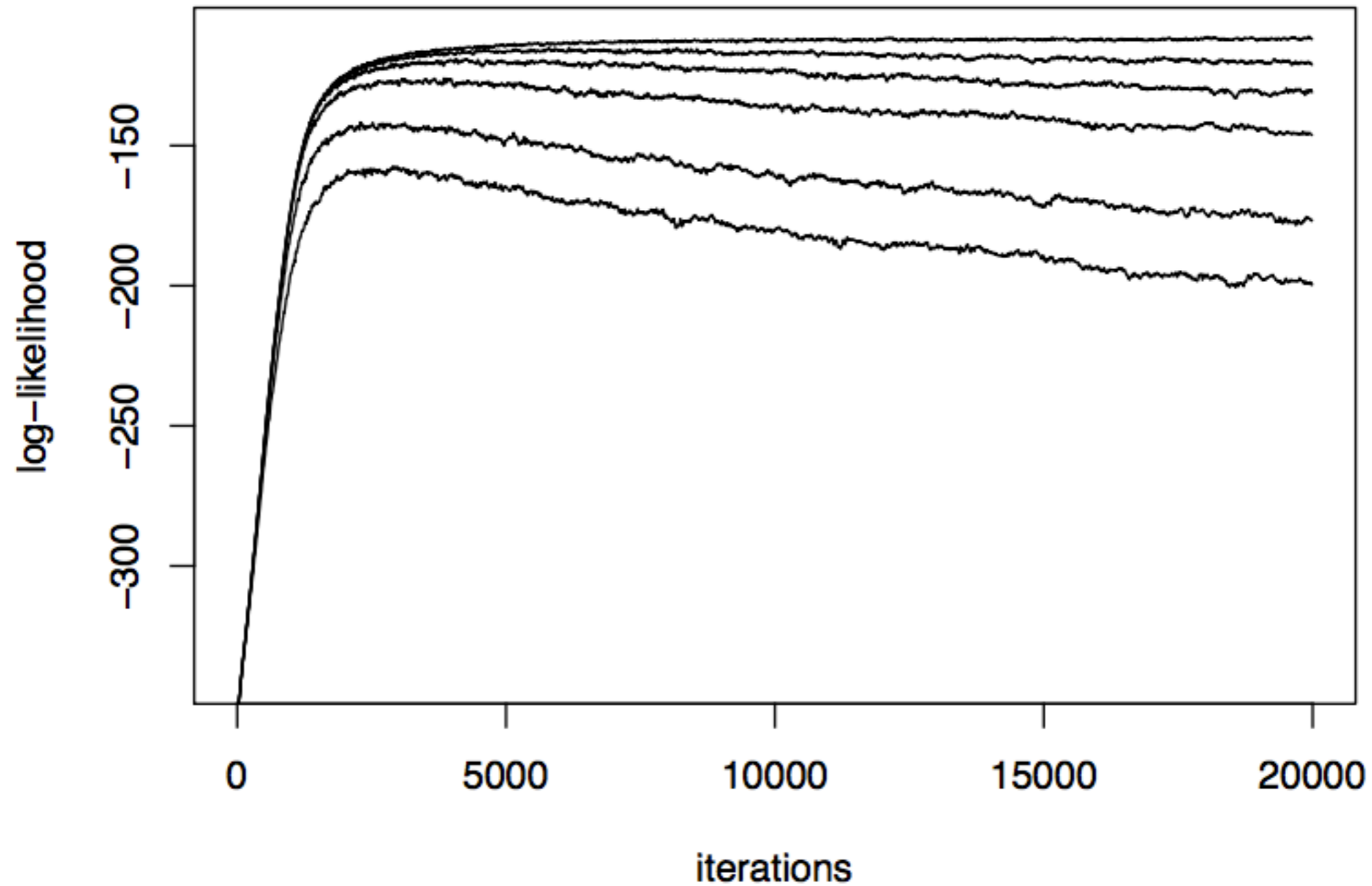
Algorithm 1: *k*-step contrastive divergence

Input: RBM $(V_1, \dots, V_m, H_1, \dots, H_n)$, training batch S

Output: gradient approximation Δw_{ij} , Δb_j and Δc_i for $i = 1, \dots, n$, $j = 1, \dots, m$

```
1 init  $\Delta w_{ij} = \Delta b_j = \Delta c_i = 0$  for  $i = 1, \dots, n$ ,  $j = 1, \dots, m$ 
2 forall the  $v \in S$  do
3    $v^{(0)} \leftarrow v$ 
4   for  $t = 0, \dots, k - 1$  do
5     for  $i = 1, \dots, n$  do sample  $h_i^{(t)} \sim p(h_i | v^{(t)})$ 
6     ;
7     for  $j = 1, \dots, m$  do sample  $v_j^{(t+1)} \sim p(v_j | h^{(t)})$ 
8     ;
9   for  $i = 1, \dots, n$ ,  $j = 1, \dots, m$  do
10     $\Delta w_{ij} \leftarrow \Delta w_{ij} + p(H_i = 1 | v^{(0)}) \cdot v_j^{(0)} - p(H_i = 1 | v^{(k)}) \cdot v_j^{(k)}$ 
11  for  $j = 1, \dots, m$  do
12     $\Delta b_j \leftarrow \Delta b_j + v_j^{(0)} - v_j^{(k)}$ 
13  for  $i = 1, \dots, n$  do
14     $\Delta c_i \leftarrow \Delta c_i + p(H_i = 1 | v^{(0)}) - p(H_i = 1 | v^{(k)})$ 
```

CD-k Running Examples



Note: from bottom to top $k = 1, 2, 5, 10, 20, 100$

Other Algorithms for Training RBM

- Persistent Contrastive Divergence (PCD)
- Fast Persistent Contrastive Divergence (FPCD)
- Stochastic Approximation Procedure (SAP)
- Stochastic Approximation Procedure with Tempered Transitions (Trans-SAP)
- Parallel Tempering (PT)

Stochastic Approximation Procedure with Tempered Transitions

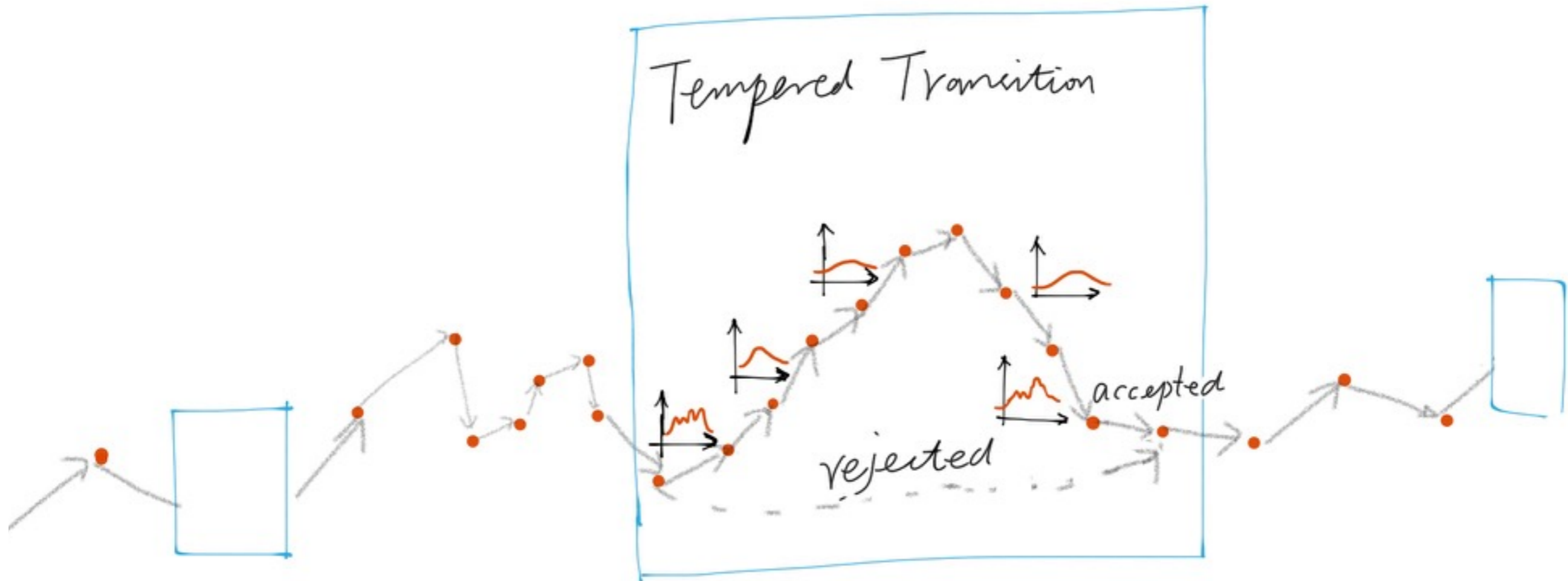
Algorithm 1 Stochastic Approximation Procedure.

- 1: Given an observation \mathbf{x}_0 . Randomly initialize θ^1 and M sample particles $\{\mathbf{x}^{1,1}, \dots, \mathbf{x}^{1,M}\}$.
 - 2: **for** $t = 1 : T$ (number of iterations) **do**
 - 3: **for** $m = 1 : M$ (number of parallel Markov chains) **do**
 - 4: Sample $\mathbf{x}^{t+1,m}$ given $\mathbf{x}^{t,m}$ using transition operator $T_{\theta^t}(\mathbf{x}^{t+1,m} \leftarrow \mathbf{x}^{t,m})$.
 - 5: **end for**
 - 6: Update: $\theta^{t+1} = \theta^t + \alpha_t \left[\phi(\mathbf{x}_0) - \frac{1}{M} \sum_{m=1}^M \phi(\mathbf{x}^{t+1,m}) \right]$.
 - 7: Decrease α_t .
 - 8: **end for**
-

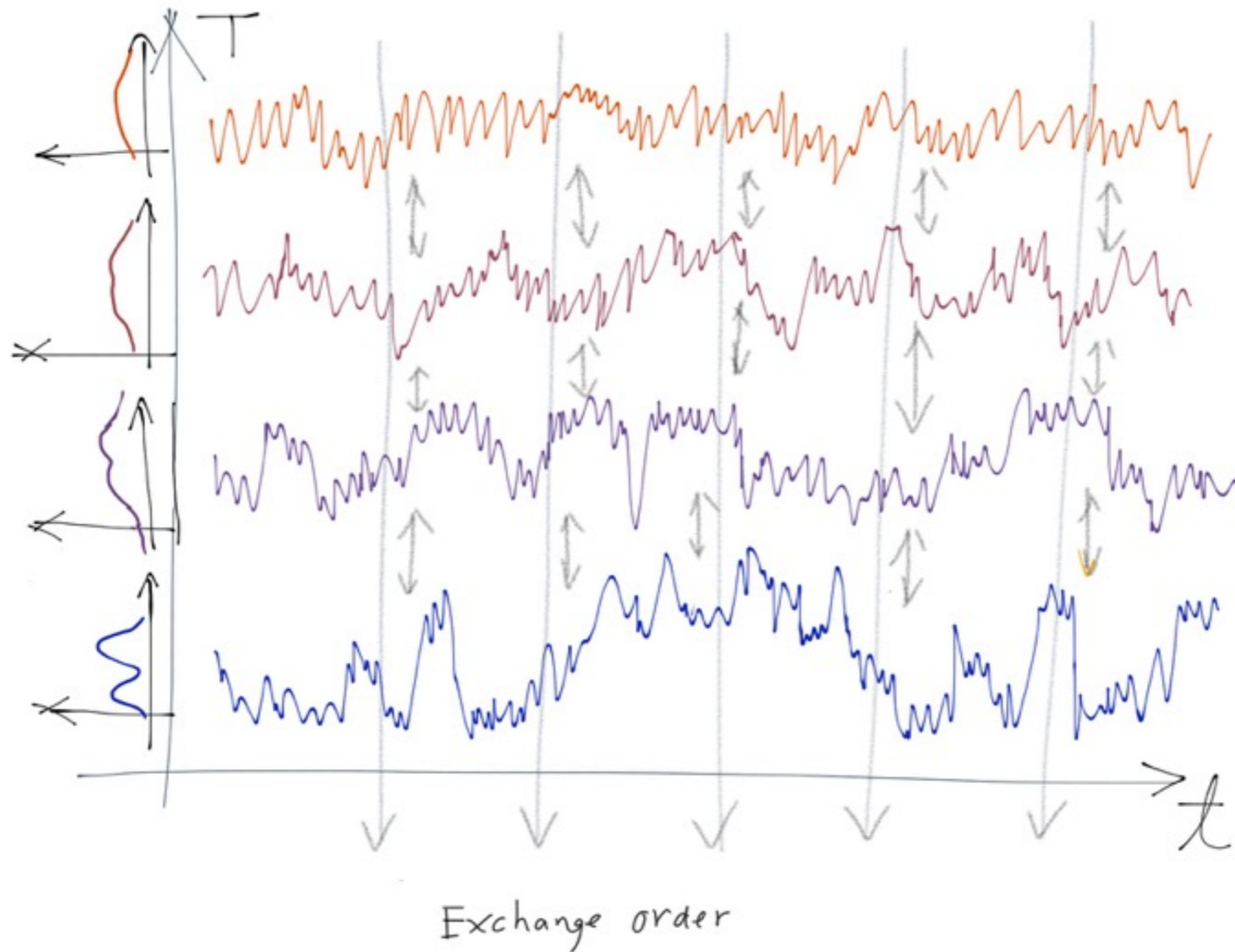
Algorithm 2 Tempered Transitions Run.

- 1: Initialize $\beta_0 < \beta_1 < \dots < \beta_S = 1$. Given a current state \mathbf{x}^S .
 - 2: **for** $s = S - 1 : 0$ (Forward pass) **do**
 - 3: Sample \mathbf{x}^s given \mathbf{x}^{s+1} using $T_s(\mathbf{x}^s \leftarrow \mathbf{x}^{s+1})$.
 - 4: **end for**
 - 5: Set $\tilde{\mathbf{x}}^0 = \mathbf{x}^0$.
 - 6: **for** $s = 0 : S - 1$ (Backward pass) **do**
 - 7: Sample $\tilde{\mathbf{x}}^{s+1}$ given $\tilde{\mathbf{x}}^s$ using $\tilde{T}_s(\tilde{\mathbf{x}}^{s+1} \leftarrow \tilde{\mathbf{x}}^s)$.
 - 8: **end for**
 - 9: Accept a new state $\tilde{\mathbf{x}}^S$ with probability: $\min \left[1, \prod_{s=1}^S p^*(\mathbf{x}_s)^{\beta_{s-1} - \beta_s} p^*(\tilde{\mathbf{x}}_s)^{\beta_s - \beta_{s-1}} \right]$.
-


Stochastic Approximation Procedure with Tempered Transitions




Parallel Tempering



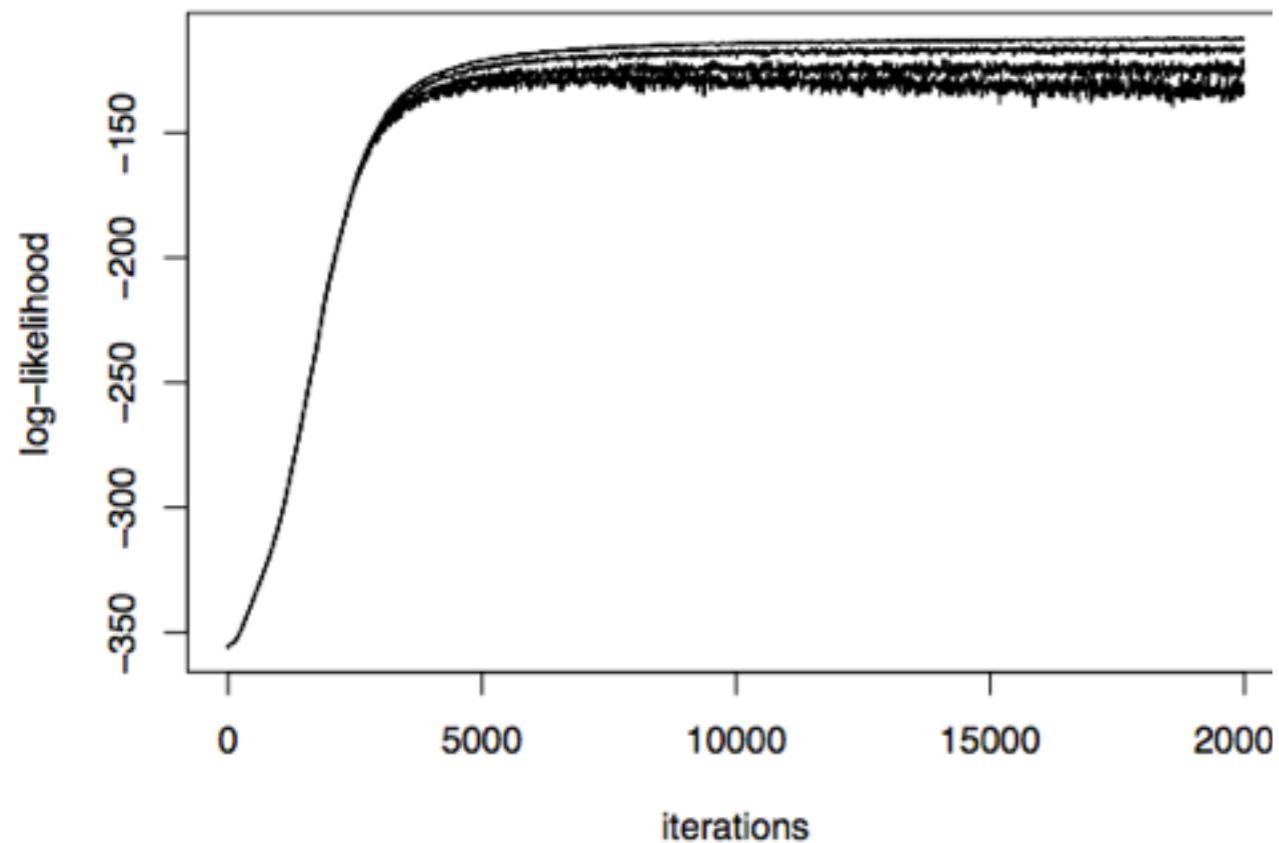
Parallel Tempering

-
- 1) Create a sequence of RBMs (R_0, R_1, \dots, R_K) such that parameters of R_k are $\Theta_k = (T_k \mathbf{W}, \mathbf{b}, \mathbf{c})$, where $0 \leq T_0 < T_1 < \dots < T_K = 1$.
 - 2) Create an empty set of samples $\mathbf{X} = \{\}$.
 - 3) Set $\mathbf{x}_0 = (\mathbf{x}_{0,0}, \dots, \mathbf{x}_{K,0})$ such that every $\mathbf{x}_{k,0}$ is a uniformly distributed random vector (or use old ones from the previous epoch).
 - 4) For $m = 1$ to M , repeat
 - a) Sample $\mathbf{x}_m = (\mathbf{x}_{0,m}, \dots, \mathbf{x}_{K,m})$ from the sequence of RBMs such that $\mathbf{x}_{k,m}$ is sampled by one-step Gibbs sampling starting from $\mathbf{x}_{k,m-1}$.
 - b) For $j = K$ to 1, repeat
 - Swap $\mathbf{x}_{j,m}$ and $\mathbf{x}_{j-1,m}$ according to $P_{\text{swap}}(\mathbf{x}_{j,m}, \mathbf{x}_{j-1,m})$ computed using .
 - c) Add $\mathbf{x}_{K,m}$ to \mathbf{X} .
 - 5) \mathbf{X} is the set of samples collected by parallel tempering sampling.
-


$$P_{\text{swap}}(\mathbf{x}_{T_1}, \mathbf{x}_{T_2}) = \min \left(1, \frac{P_{T_1}(\mathbf{x}_{T_2})P_{T_2}(\mathbf{x}_{T_1})}{P_{T_1}(\mathbf{x}_{T_1})P_{T_2}(\mathbf{x}_{T_2})} \right),$$

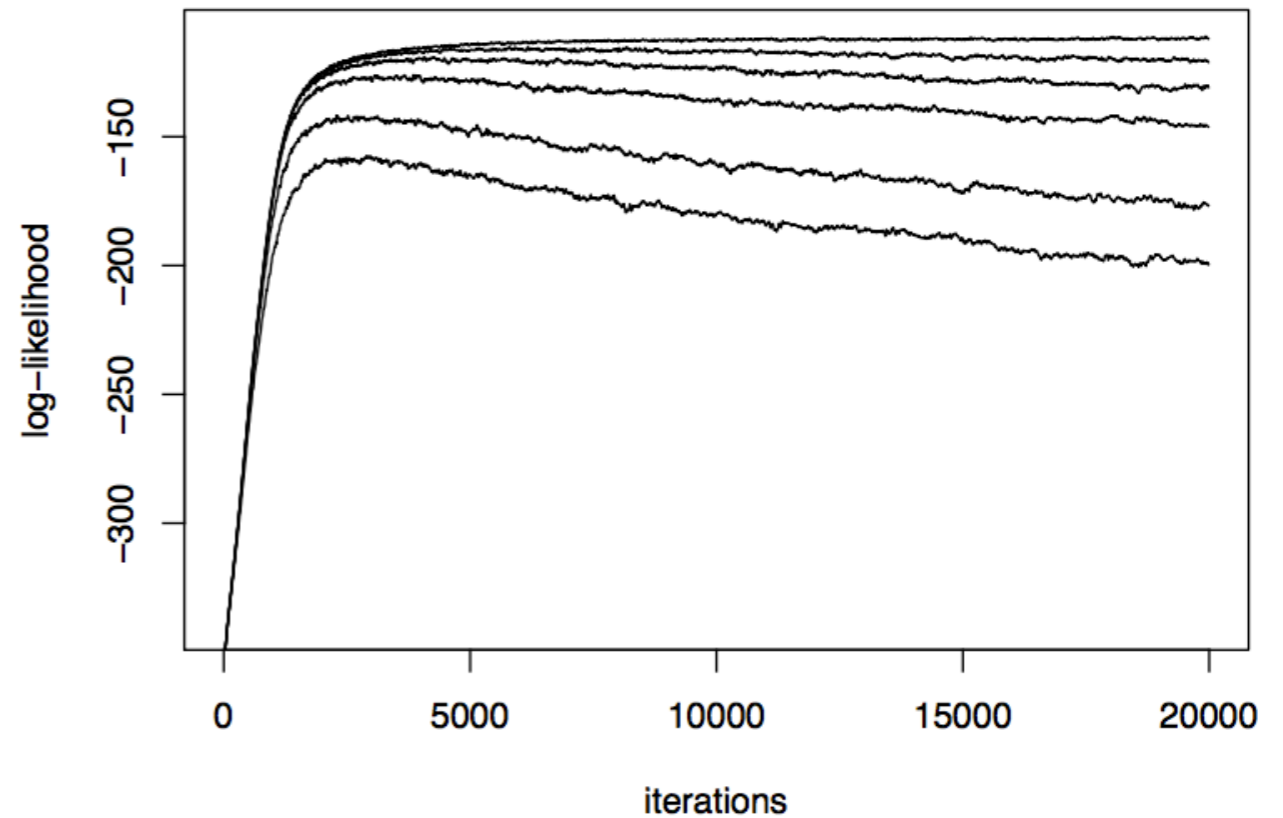
Parallel Tempering v.s. CD-k: Running Examples (1/2)

PT



Note: bottom to top $M = 4, 5, 10, 50$

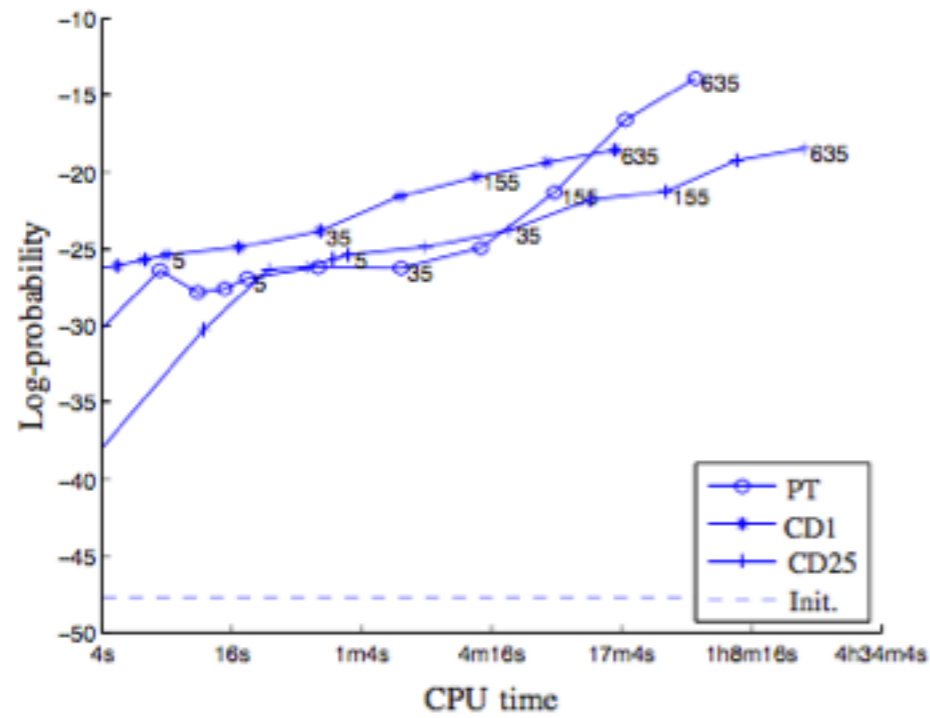
CD-k



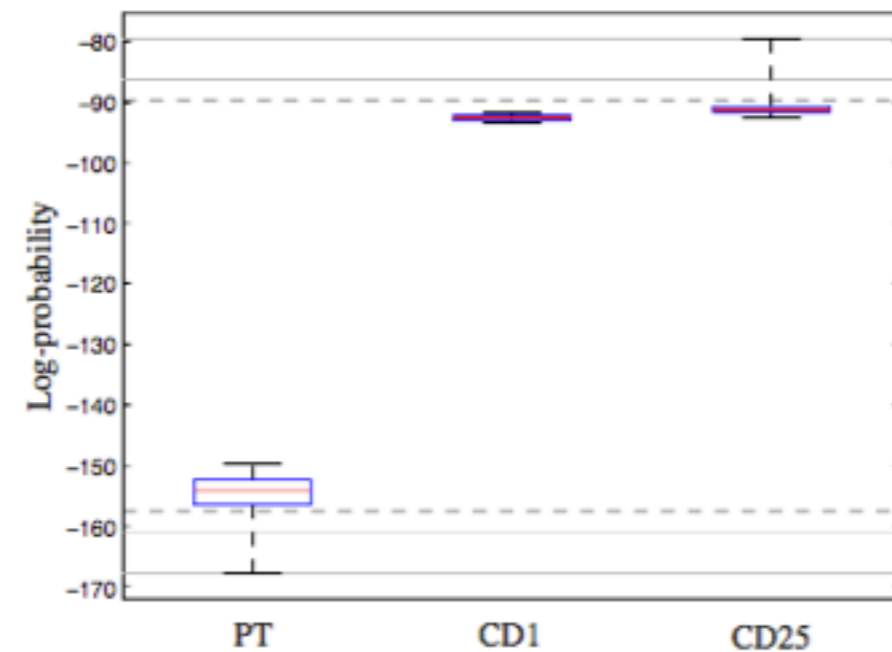
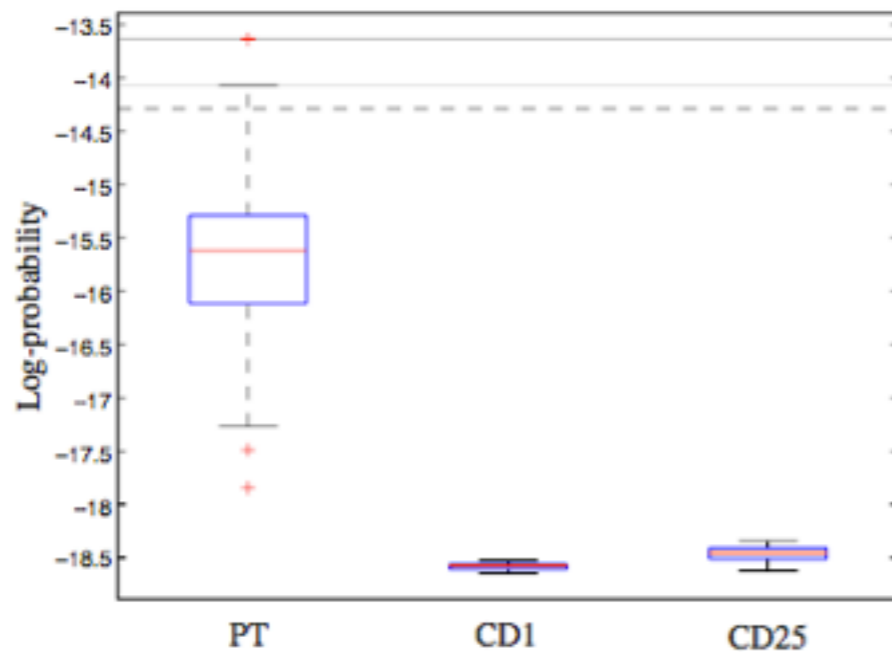
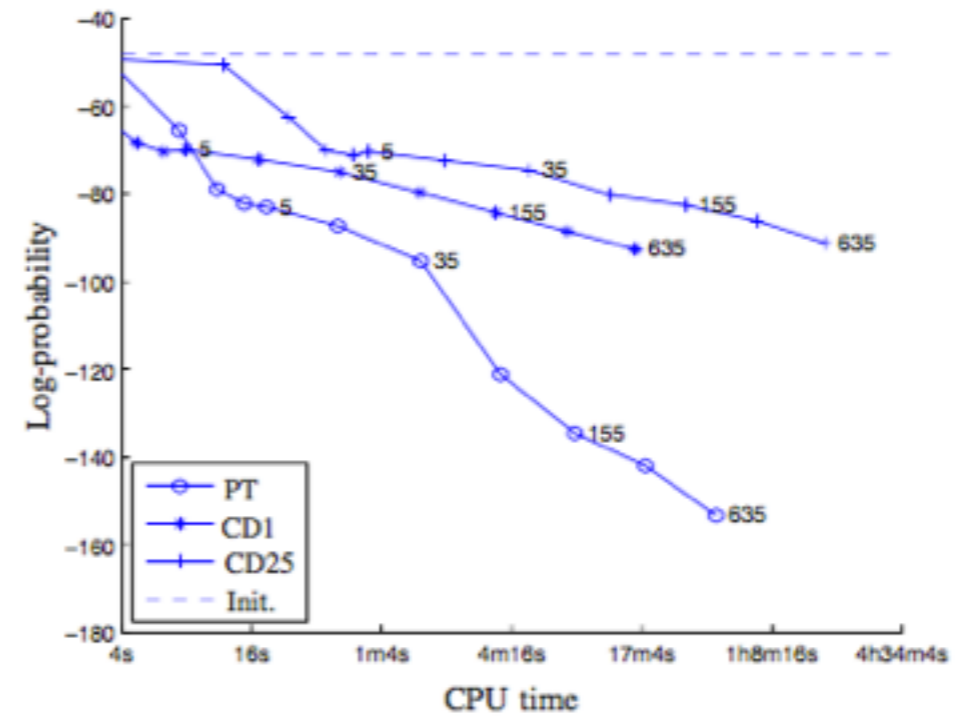
Note: from bottom to top $k = 1, 2, 5, 10, 20, 100$

Parallel Tempering v.s. CD-k: Running Examples (2/2)

Positive Samples



Negative Samples



References

- Andrieu, Christophe, Nando de Freitas, Arnaud Doucet, and Michael I Jordan. 'An Introduction to MCMC for Machine Learning.', *Machine learning*, 2003
- Fischer, Asja, and Christian Igel. 'Training restricted Boltzmann machines: An introduction.', *Pattern Recognition*, 2014
- Salakhutdinov, Ruslan. 'Learning in Markov Random Fields using Tempered Transitions.', NIPS, 2009
- Geoffrey E. Hinton, Simon Osindero, Simon Osindero, and Yee-Whye Teh. 'A Fast Learning Algorithm for Deep Belief Nets', *Neural computation*, 2006
- Desjardins, Guillaume, Aaron C Courville, Yoshua Bengio, Pascal Vincent, and Olivier Delalleau. 'Parallel Tempering for Training of Restricted Boltzmann Machines ', *AISTATS*, 2010
- Neal, Radford M. 'Sampling from multimodal distributions using tempered transitions', *Statistics and computing*, 1996
- Geoffrey E. Hinton. 'Reducing the Dimensionality of Data with Neural Networks', *Science*, 2006
- Cho, Kyunghyun, Tapani Raiko, and Alexander Ilin. 'Parallel tempering is efficient for learning restricted Boltzmann machines', *International Joint Conference on Neural Networks (IJCNN)*, 2010

Thanks