

AUSTERITY IN MCMC LAND: CUTTING THE METROPOLIS-HASTINGS BUDGET

Anoop Korattikara



Yutian Chen



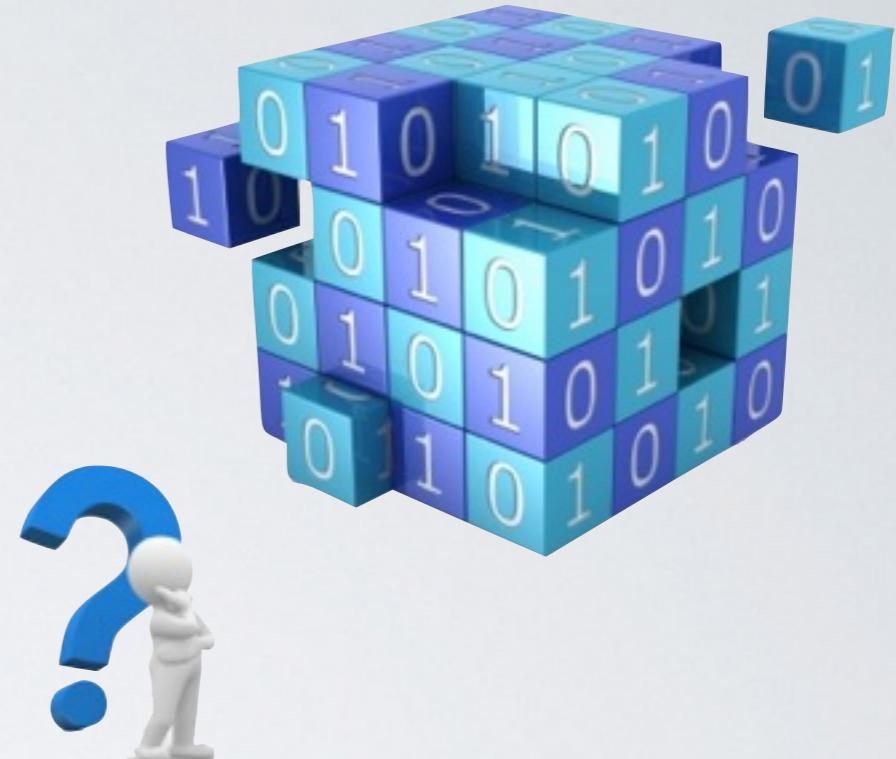
Max Welling



UNIVERSITEIT VAN AMSTERDAM

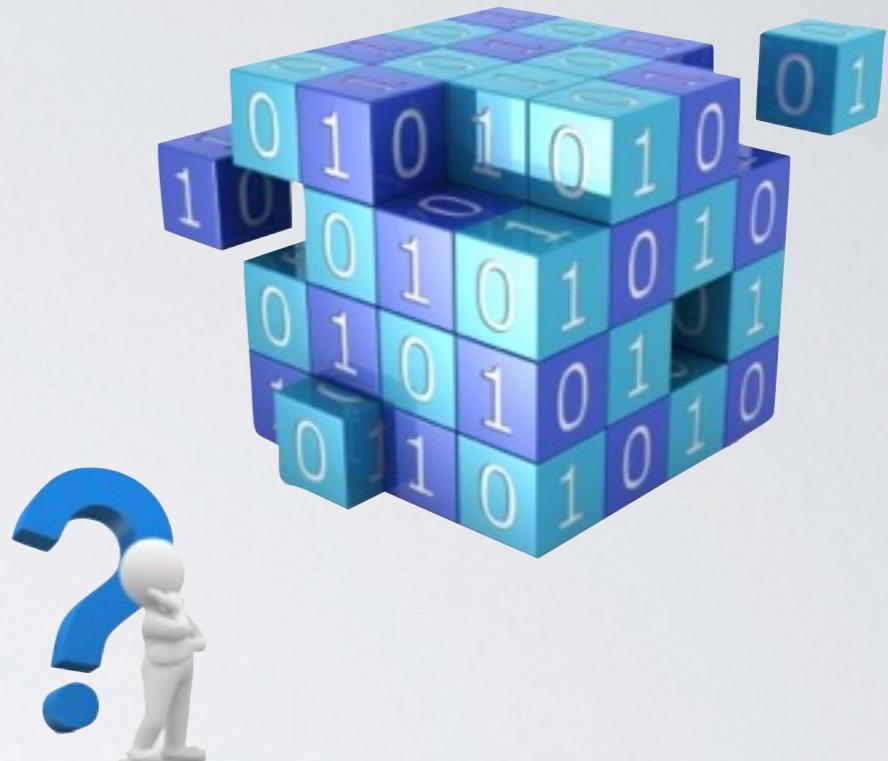
University of Sheffield, 6 March 2014

Big Data Challenge for Machine Learning



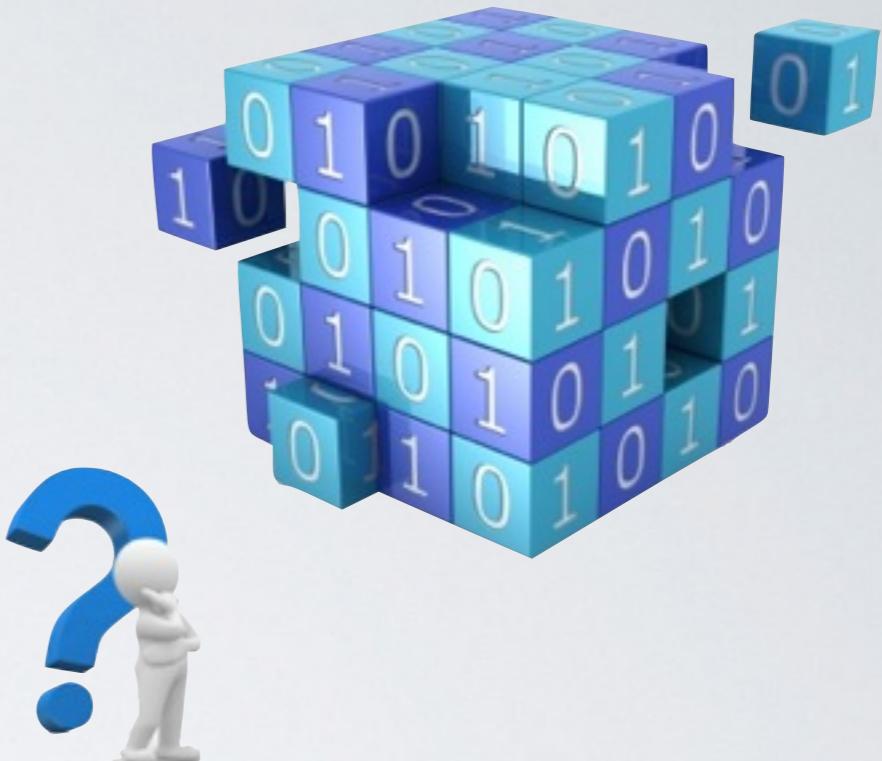
Big Data Challenge for Machine Learning

- Standard settings, moderate size of data
 - Generalization error, $E(N)$
 - Computational efficiency, $T(N)$

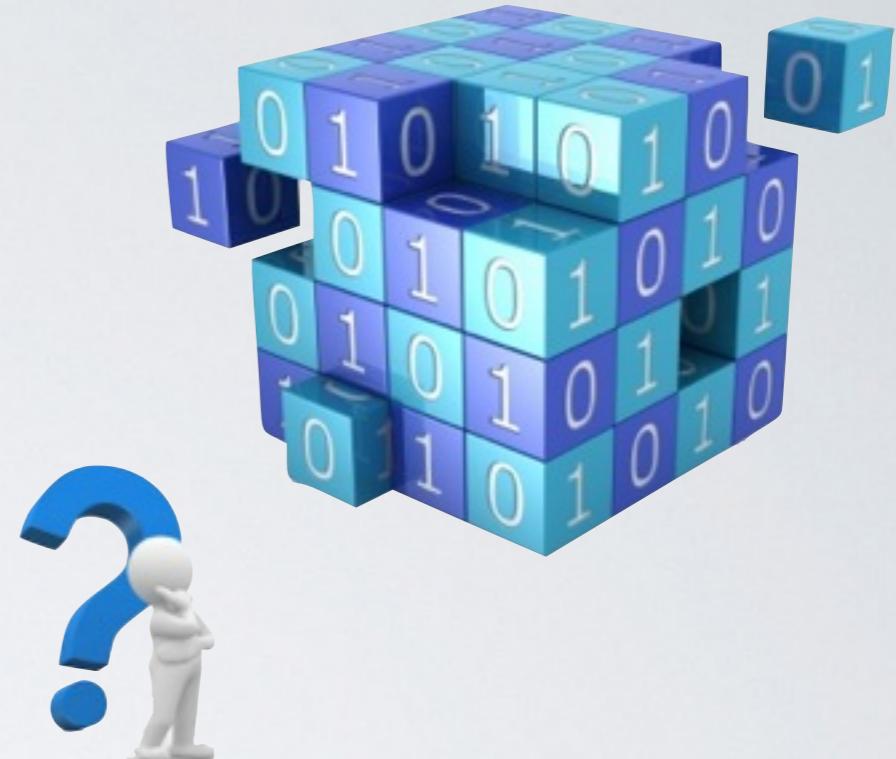


Big Data Challenge for Machine Learning

- Standard settings, moderate size of data
 - Generalization error, $E(N)$
 - Computational efficiency, $T(N)$
- Big data settings
 - Limited budget on time
 - $E(N, T)$ or $E(T)$ (Bottou & Bousquet, 2007)



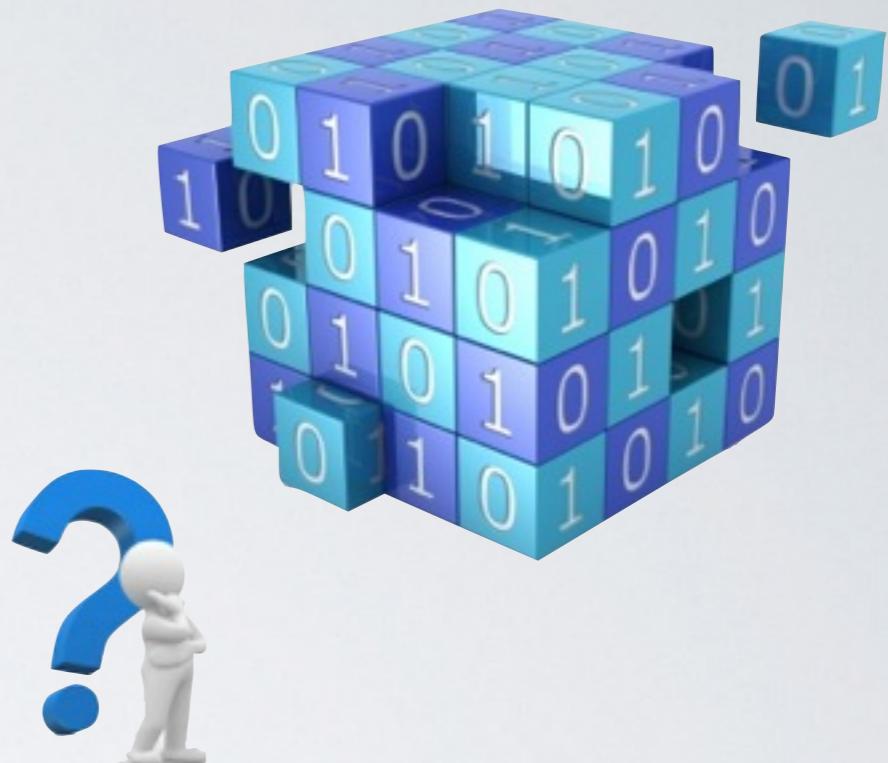
Big Data Challenge for Machine Learning



Big Data Challenge for Machine Learning

- Learning with optimization

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N l(x_i; \theta)$$

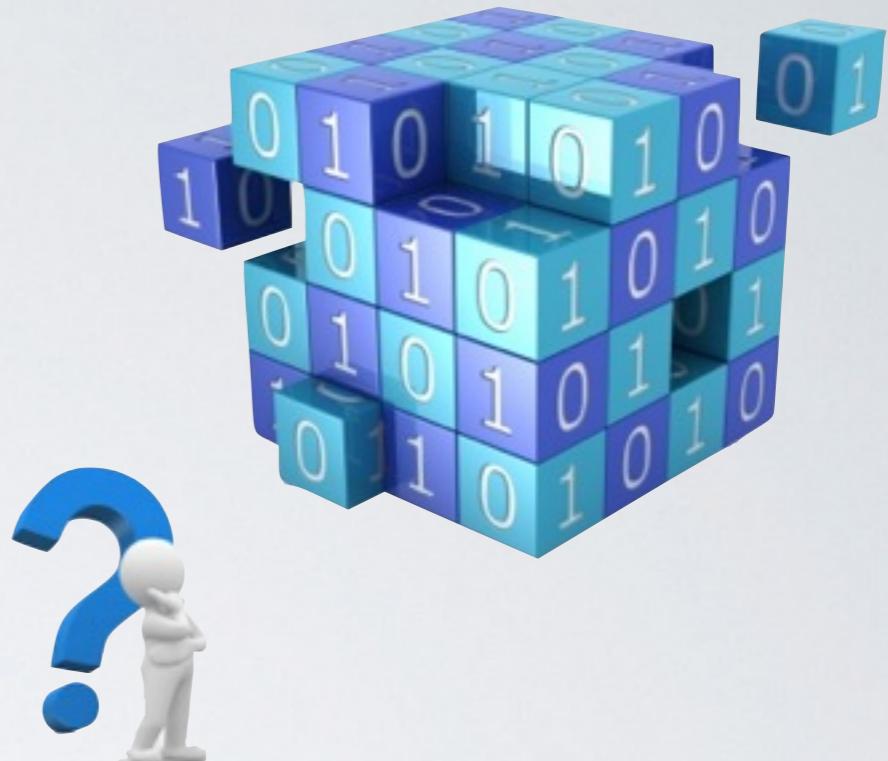


Big Data Challenge for Machine Learning

- Learning with optimization

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N l(x_i; \theta)$$

O(N)



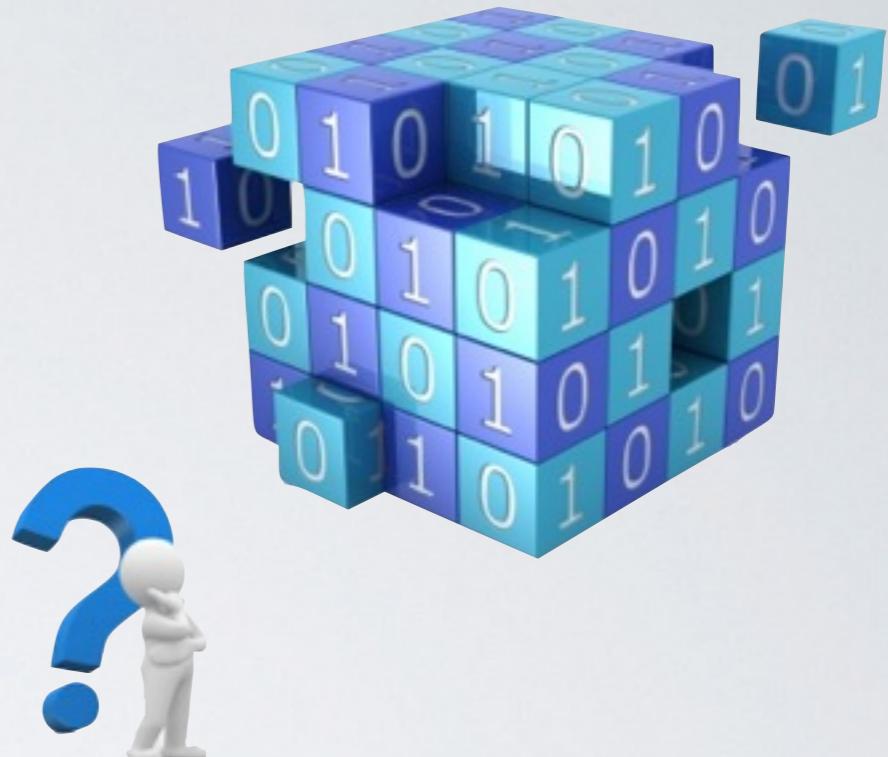
Big Data Challenge for Machine Learning

- Learning with optimization

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N l(x_i; \theta)$$

O(N)

- Stochastic gradient descent **O(n)**



Big Data Challenge for Machine Learning

- Learning with optimization

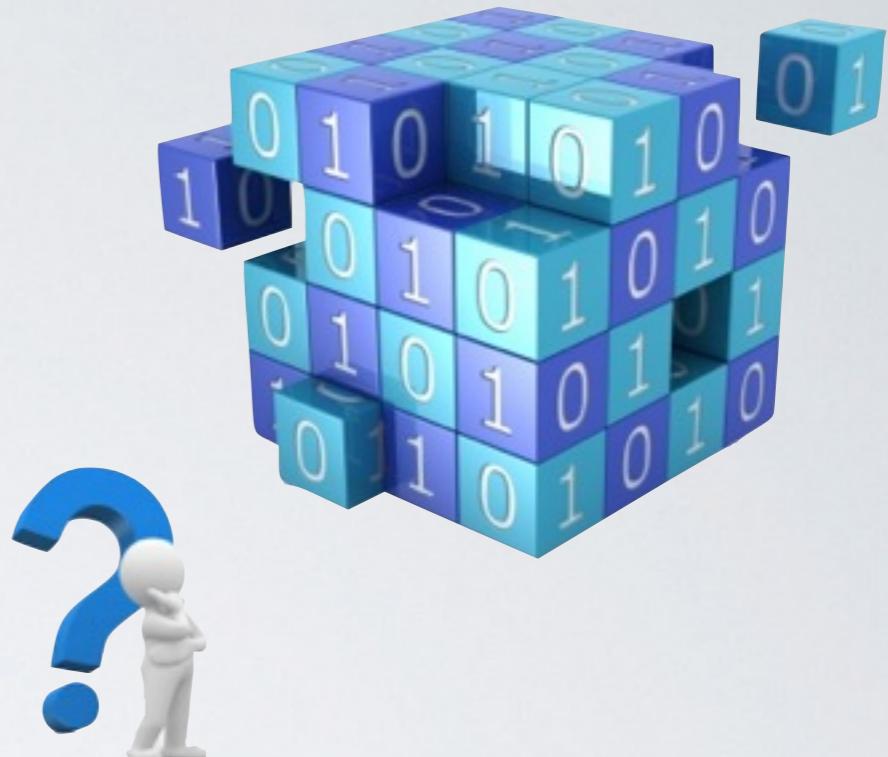
$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N l(x_i; \theta)$$

O(N)

- Stochastic gradient descent **O(n)**

- Bayesian inference

$$P(\theta | \mathcal{D}) \propto P(\theta) \prod_{i=1}^N P(x_i | \theta)$$



Big Data Challenge for Machine Learning

- Learning with optimization

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N l(x_i; \theta)$$

O(N)

- Stochastic gradient descent **O(n)**

- Bayesian inference

$$P(\theta | \mathcal{D}) \propto P(\theta) \prod_{i=1}^N P(x_i | \theta) \quad \textcolor{red}{O(N)}$$



Big Data Challenge for Machine Learning

- Learning with optimization

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N l(x_i; \theta)$$

O(N)

- Stochastic gradient descent **O(n)**



- Bayesian inference

$$P(\theta | \mathcal{D}) \propto P(\theta) \prod_{i=1}^N P(x_i | \theta) \quad \textcolor{red}{O(N)}$$

- Stochastic Variational Inference (Hoffman, et al, 2013) **O(n)**



Big Data Challenge for Machine Learning

- Learning with optimization

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N l(x_i; \theta)$$

O(N)

- Stochastic gradient descent **O(n)**



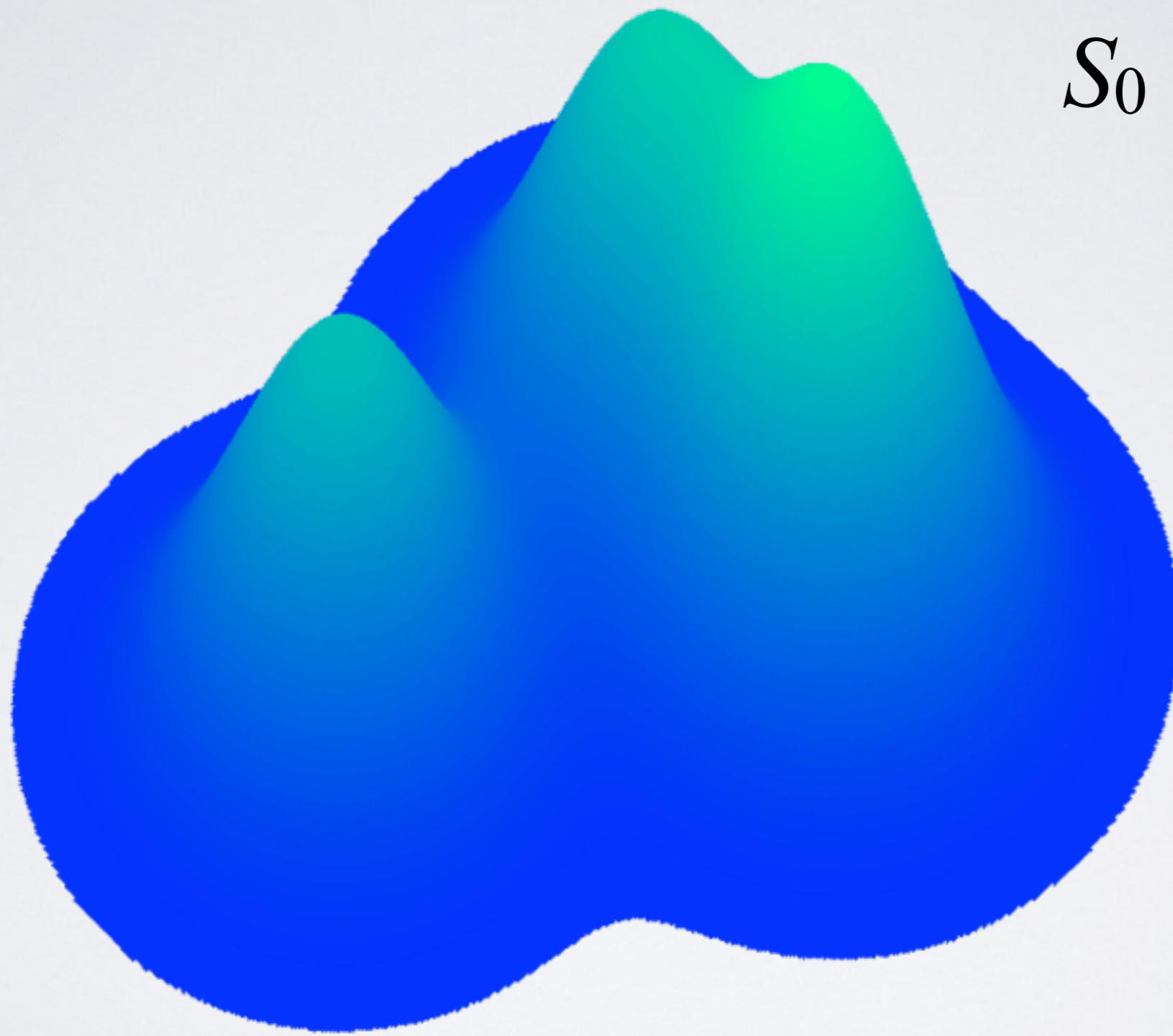
- Bayesian inference

$$P(\theta | \mathcal{D}) \propto P(\theta) \prod_{i=1}^N P(x_i | \theta) \quad \textcolor{red}{O(N)}$$

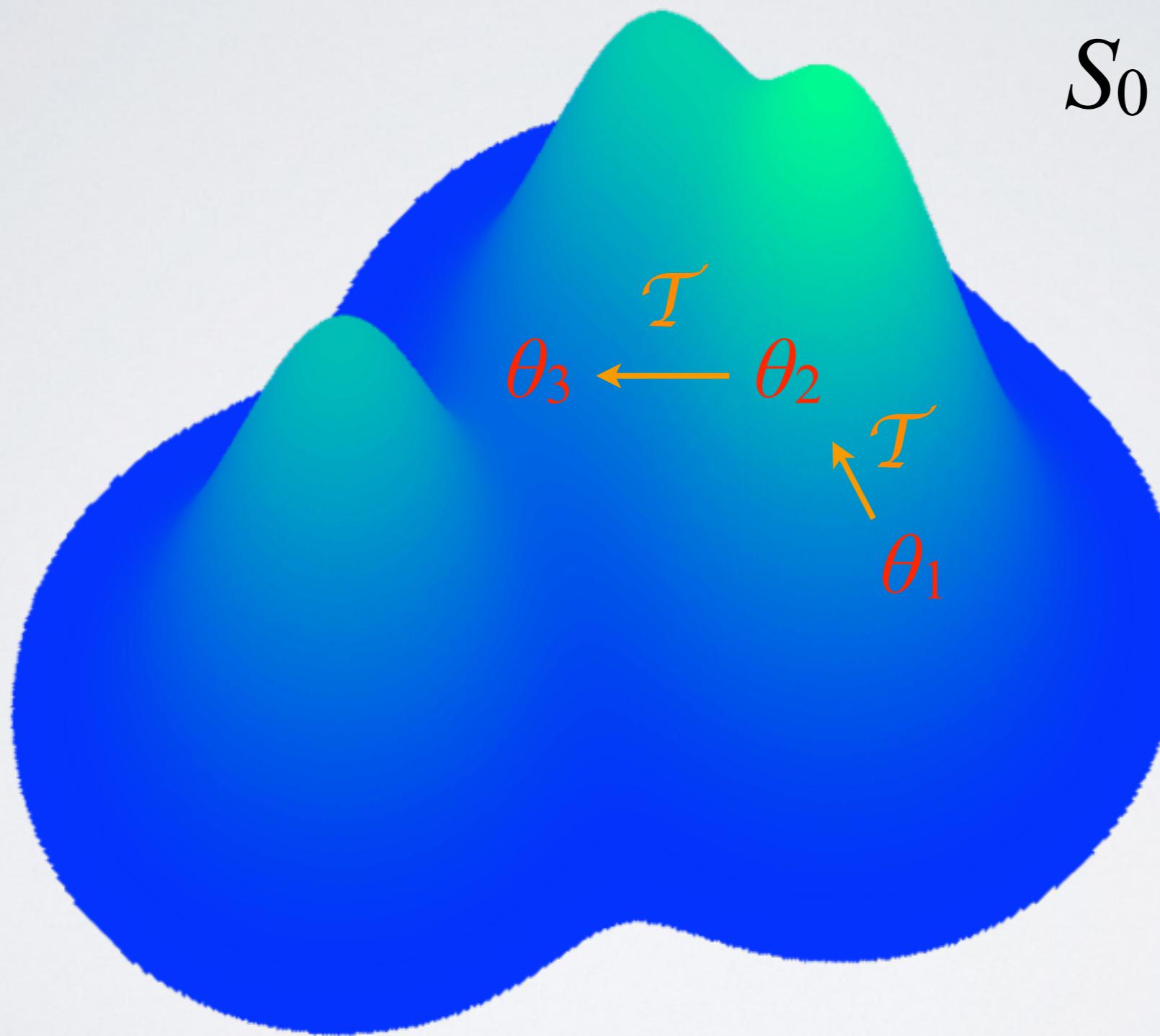
- Stochastic Variational Inference (Hoffman, et al, 2013) **O(n)**
- MCMC



MCMC Review

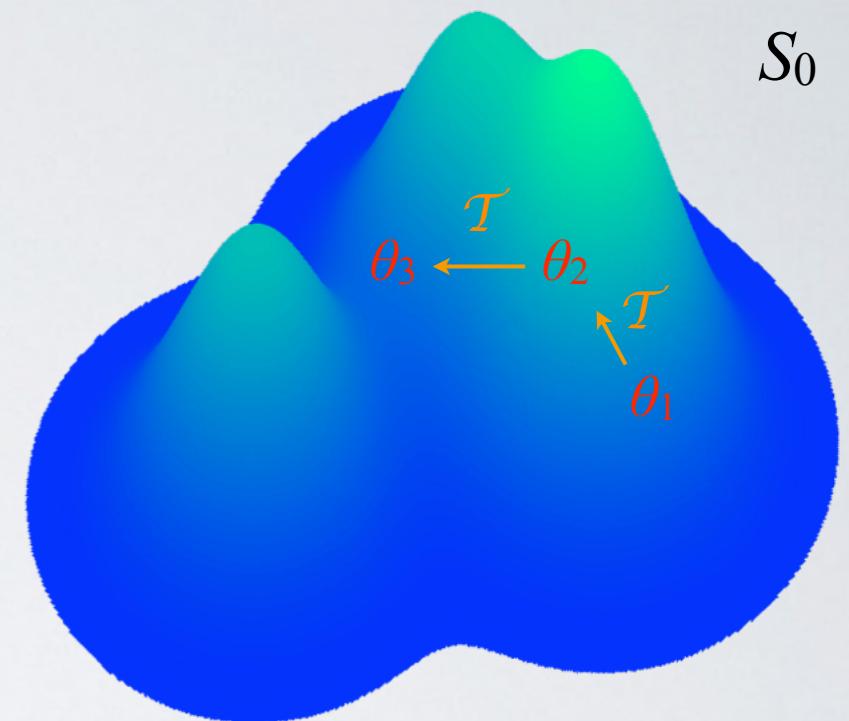


MCMC Review



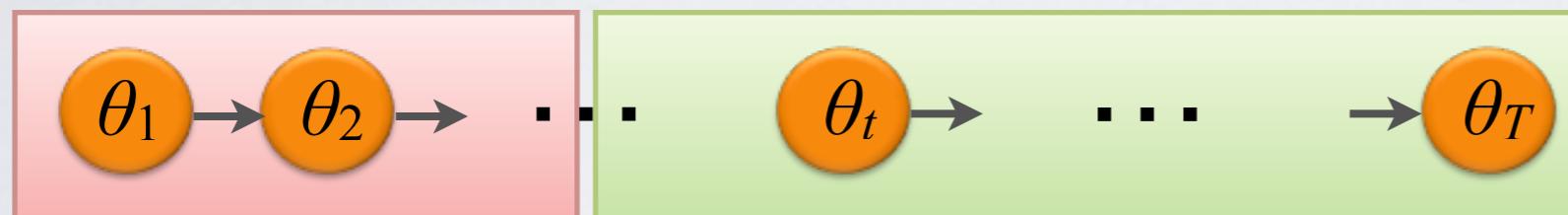
MCMC Review

Choose \mathcal{T} s.t. the marginals $P(\theta_t) \rightarrow S_0$ as $t \rightarrow \infty$

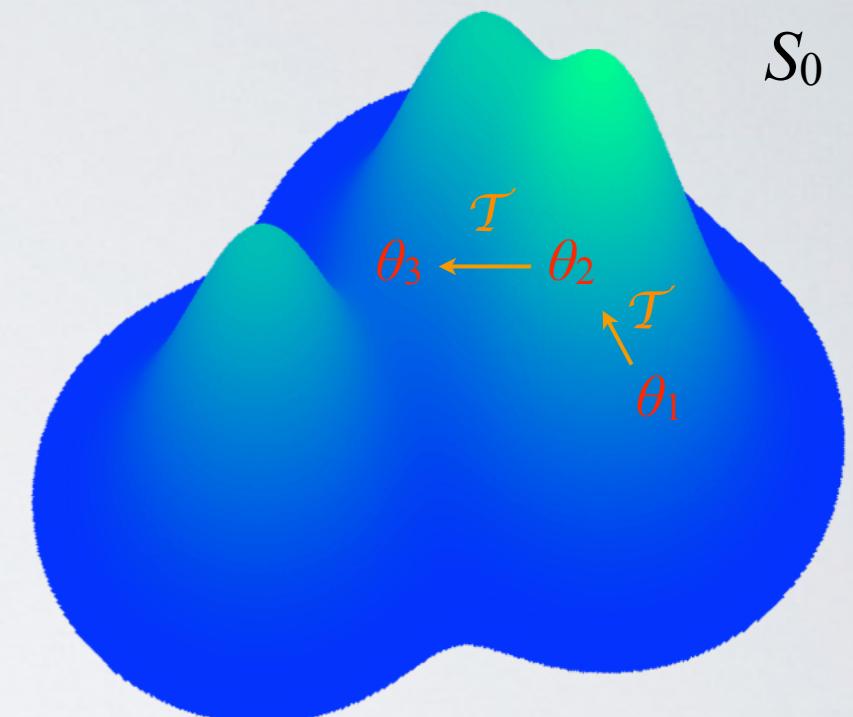


MCMC Review

Choose \mathcal{T} s.t. the marginals $P(\theta_t) \rightarrow S_0$ as $t \rightarrow \infty$

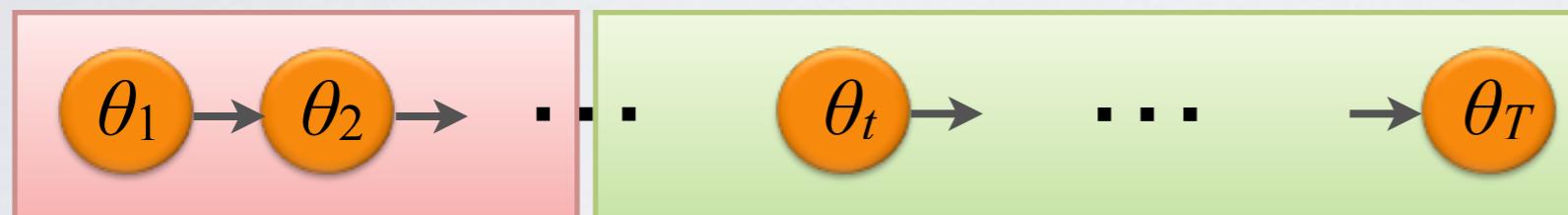


Burn-in (Throw away)



MCMC Review

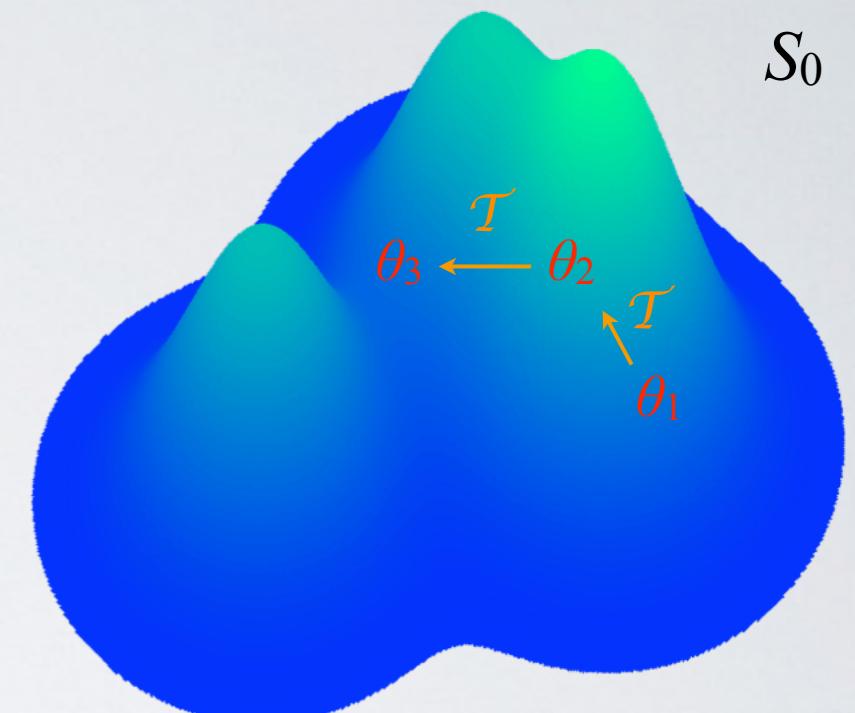
Choose \mathcal{T} s.t. the marginals $P(\theta_t) \rightarrow S_0$ as $t \rightarrow \infty$



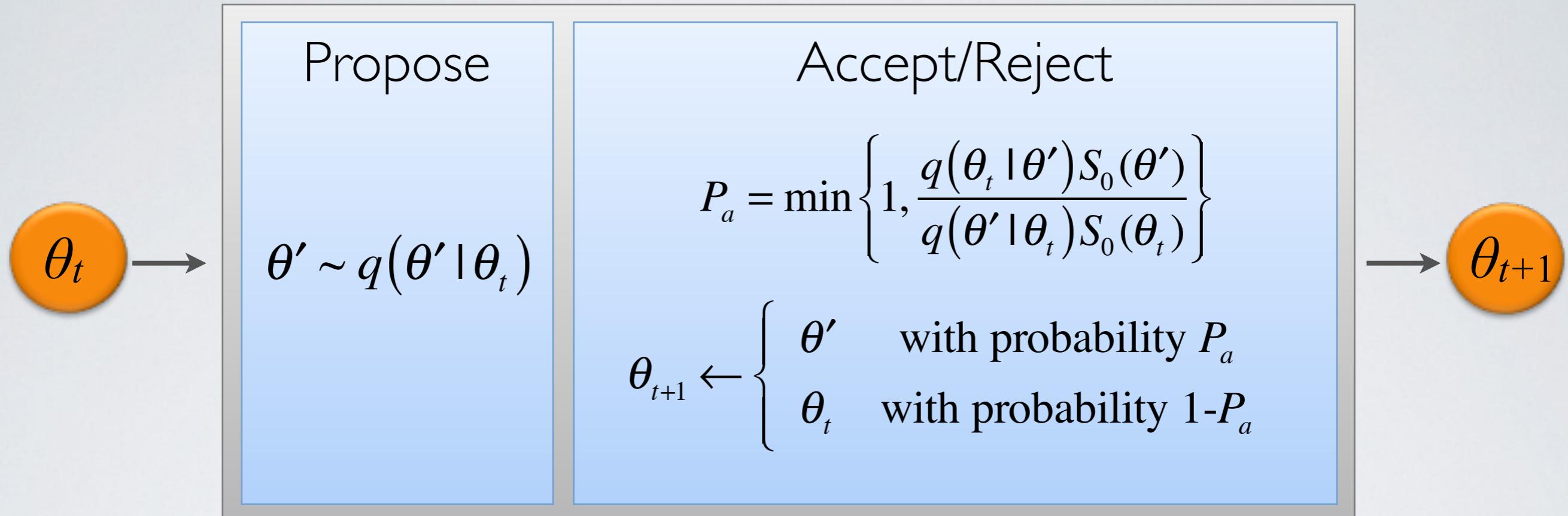
Burn-in (Throw away)

Use here

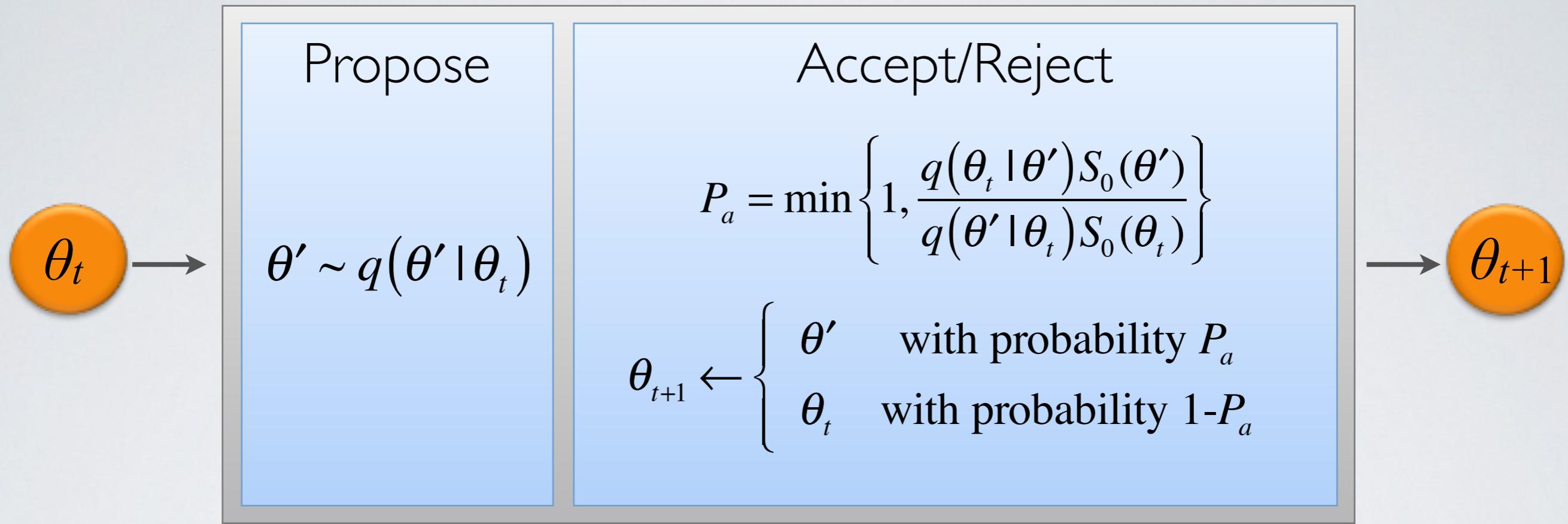
$$I = \langle f \rangle_{S_0} \approx \hat{I} = \overline{f(\theta_t)}, \text{ where } \theta_t \sim S_0$$



Metropolis-Hastings Sampling

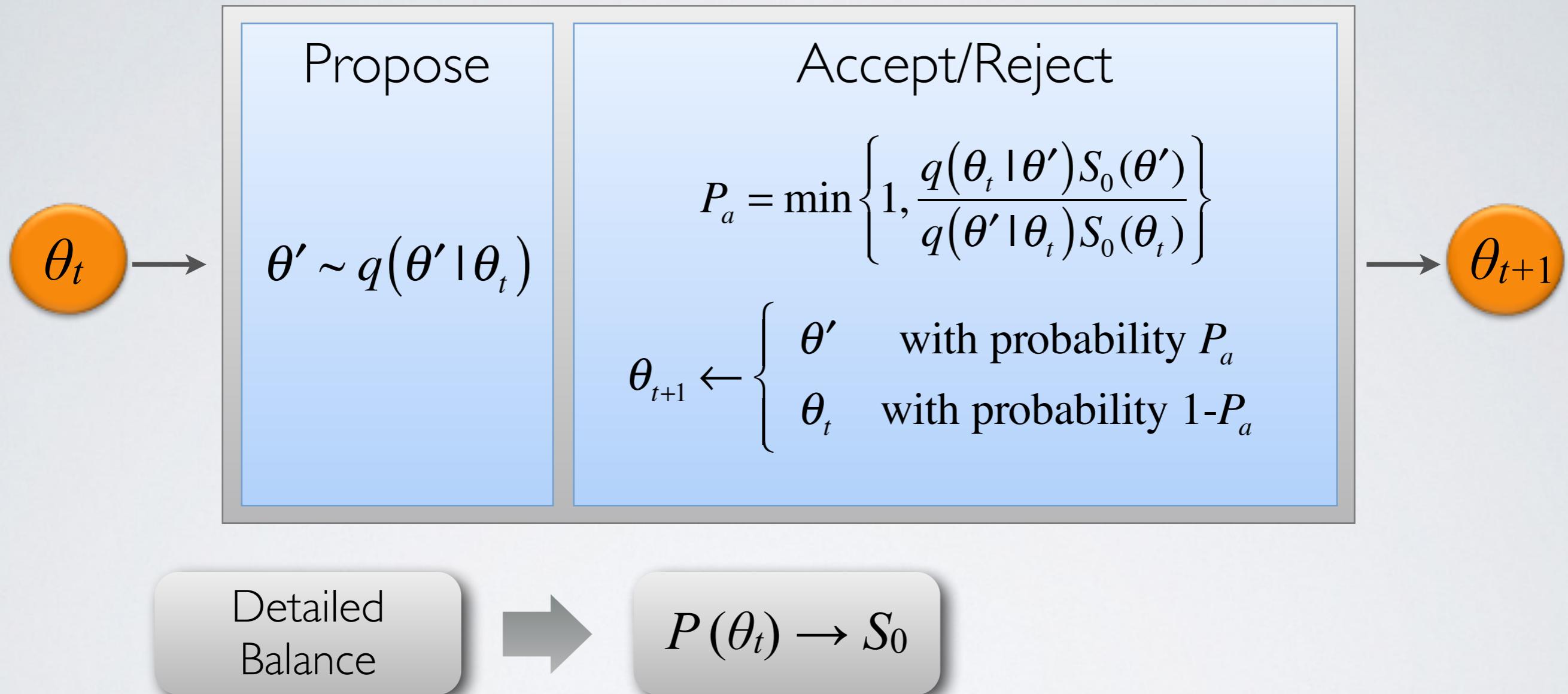


Metropolis-Hastings Sampling

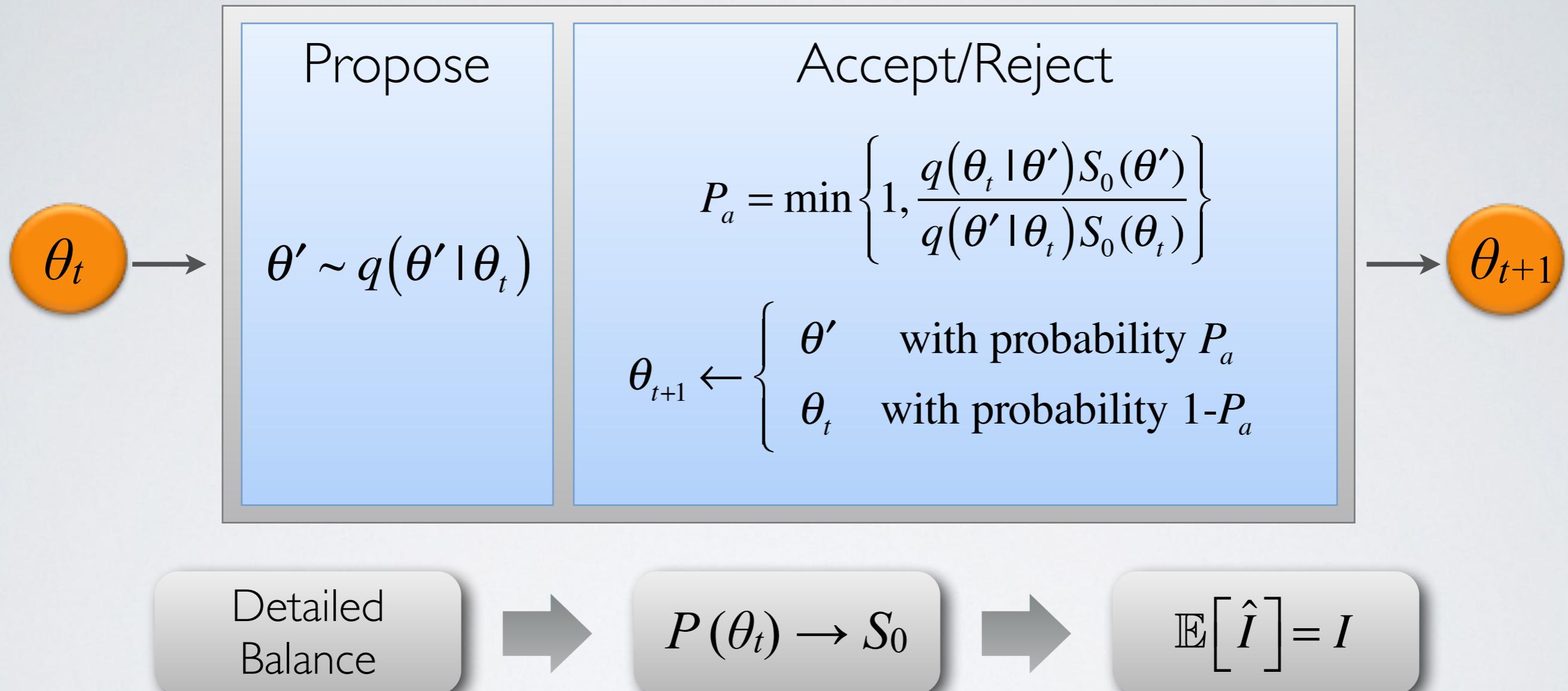


Detailed
Balance

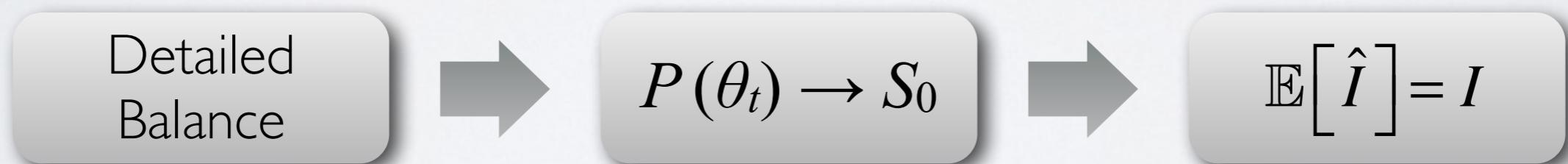
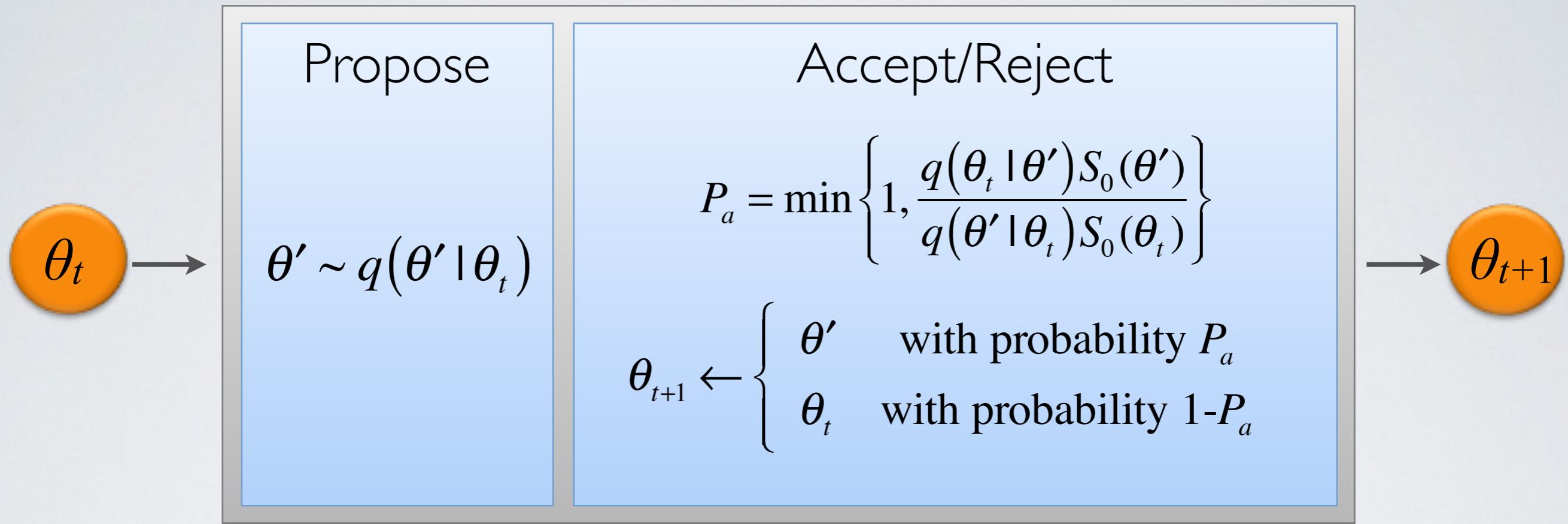
Metropolis-Hastings Sampling



Metropolis-Hastings Sampling

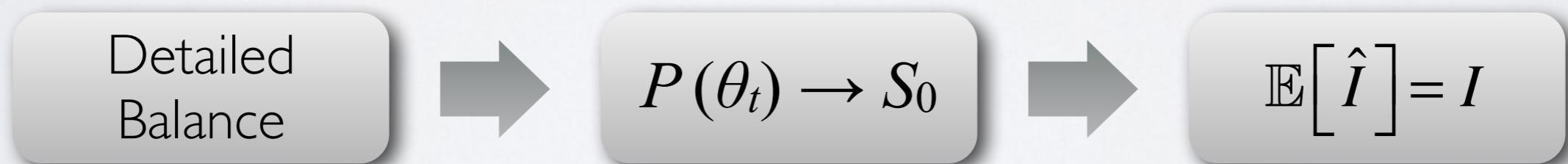
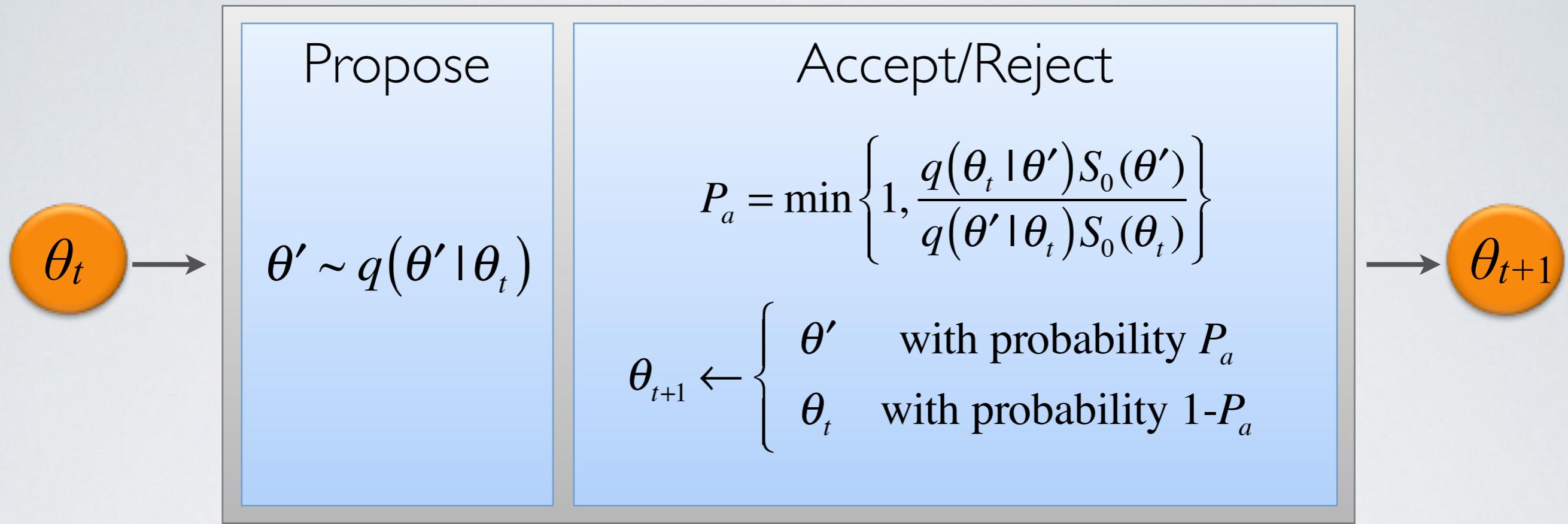


Metropolis-Hastings Sampling



$$S_0(\theta) = P(\theta | \mathcal{D}) \propto P(\theta) \prod_{i=1}^N P(x_i | \theta)$$

Metropolis-Hastings Sampling



$$S_0(\theta) = P(\theta | \mathcal{D}) \propto P(\theta) \prod_{i=1}^N P(x_i | \theta)$$

O(N)

Risk Decomposition

$$\text{Risk} = \text{Bias}^2 + \text{Variance}$$
$$\mathbb{E}\left[\left(I - \hat{I}\right)^2\right] = I - \langle f \rangle_S + \sigma^2 \frac{\tau}{T}$$

Risk Decomposition

$$\text{Risk} = \text{Bias}^2 + \text{Variance}$$

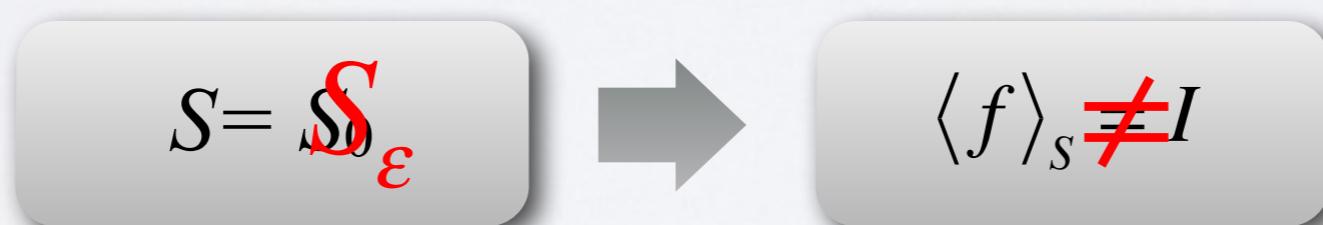
| | | | |
|---|---|---------------------------|---------------------------|
| Risk | = | Bias | Bias^2 |
| $\mathbb{E}\left[\left(I - \hat{I}\right)^2\right]$ | | $I - \langle f \rangle_S$ | $\sigma^2 \frac{\tau}{T}$ |

$S = S_0 \rightarrow \langle f \rangle_S = I$

Risk Decomposition

$$\text{Risk} = \text{Bias}^2 + \text{Variance}$$

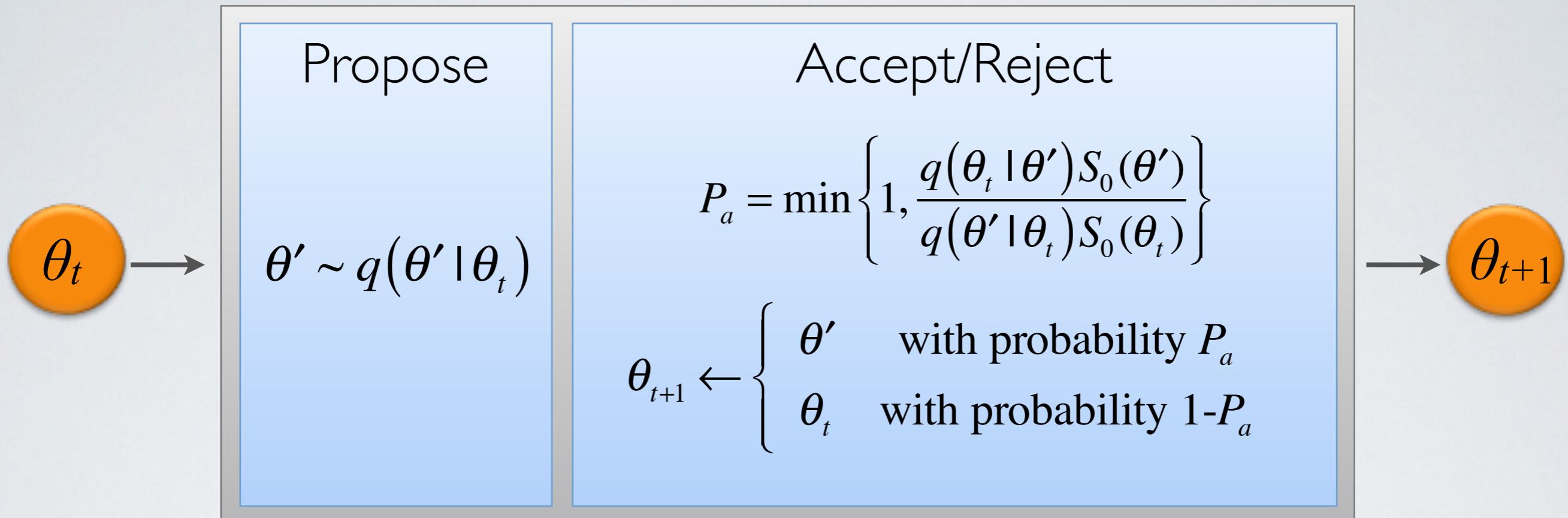
| | | | |
|---|---|---------------------------|---------------------------|
| Risk | = | Bias | Variance |
| $\mathbb{E}\left[\left(I - \hat{I}\right)^2\right]$ | | $I - \langle f \rangle_S$ | $\sigma^2 \frac{\tau}{T}$ |



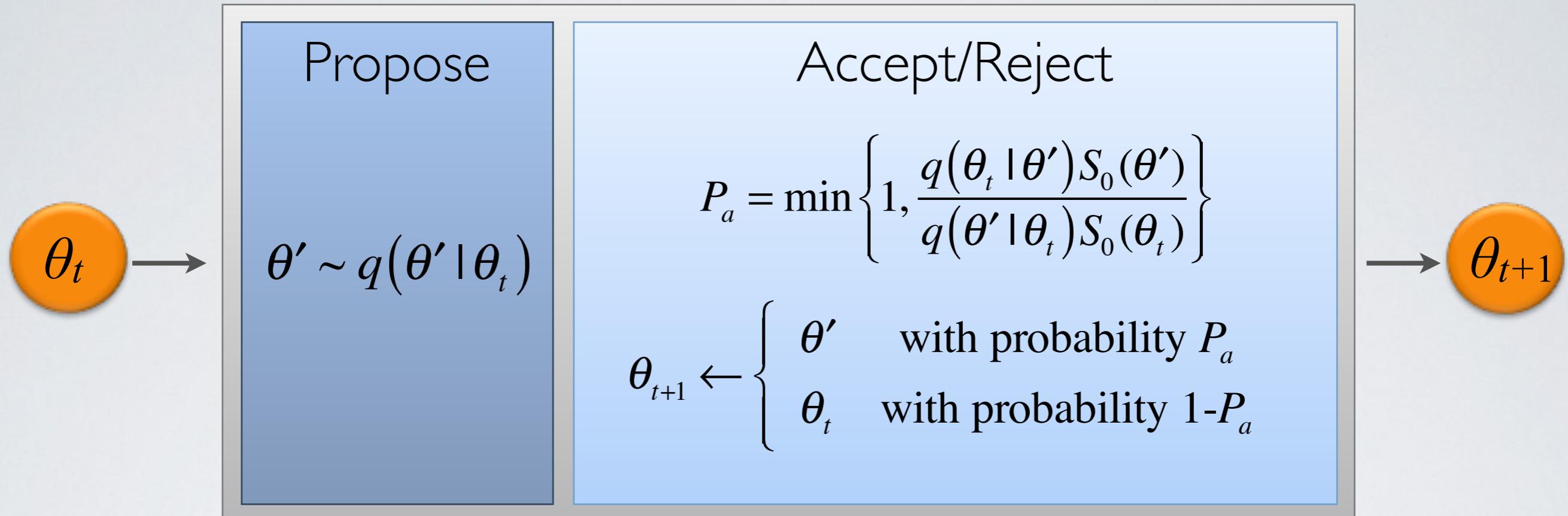
Fixed Computational Budget



Efficient Approximate MH

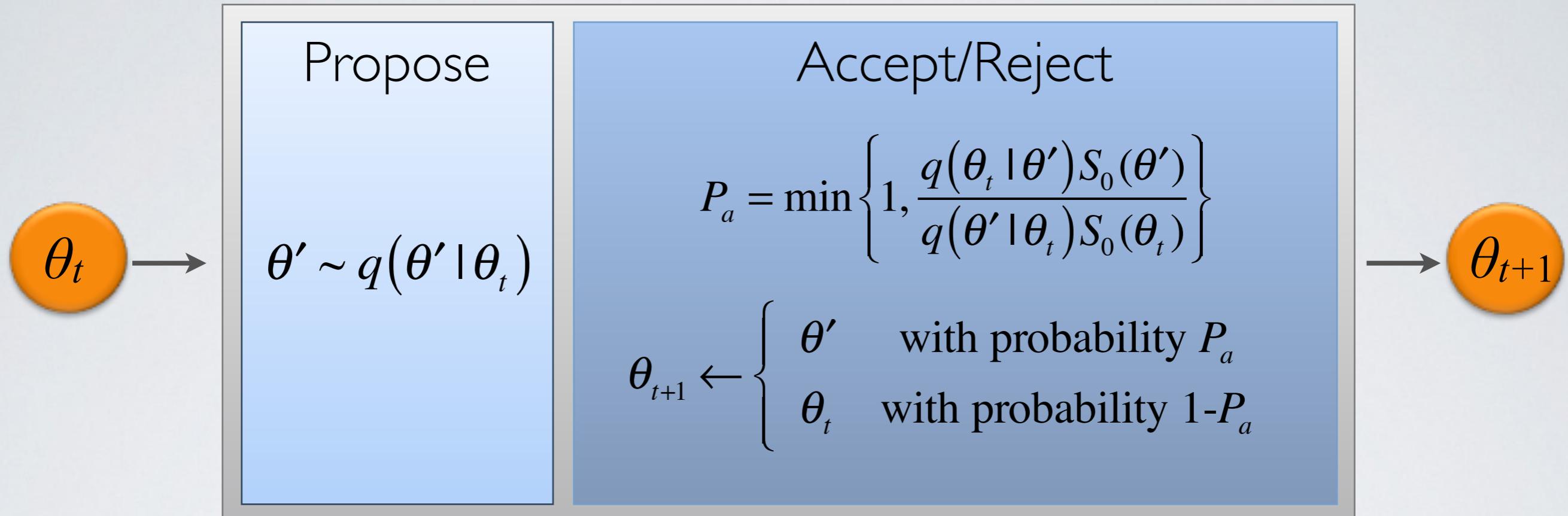


Efficient Approximate MH



- Efficient Proposal
- Stochastic Gradient Langevin Dynamics
(Welling&Teh, 2011; Ahn et.al. 2012)

Efficient Approximate MH



- Efficient Accept/Reject?

Accept/Reject as a Hypothesis Testing

$$P_a = \min \left\{ 1, \frac{q(\theta_t | \theta') S_0(\theta')}{q(\theta' | \theta_t) S_0(\theta_t)} \right\}$$

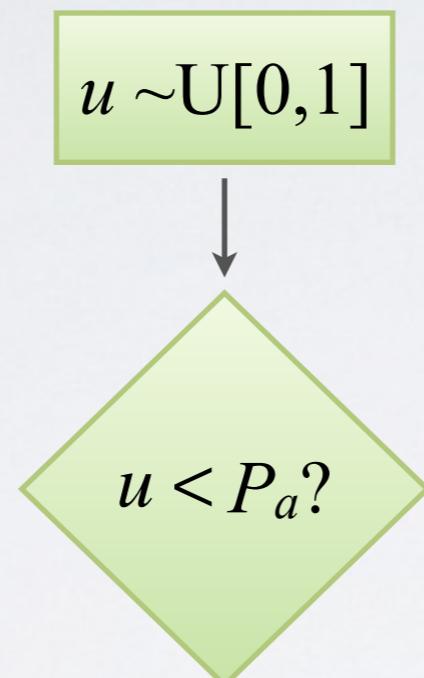
Accept/Reject as a Hypothesis Testing

$$P_a = \min \left\{ 1, \frac{q(\theta_t | \theta') S_0(\theta')}{q(\theta' | \theta_t) S_0(\theta_t)} \right\}$$

$$u \sim U[0,1]$$

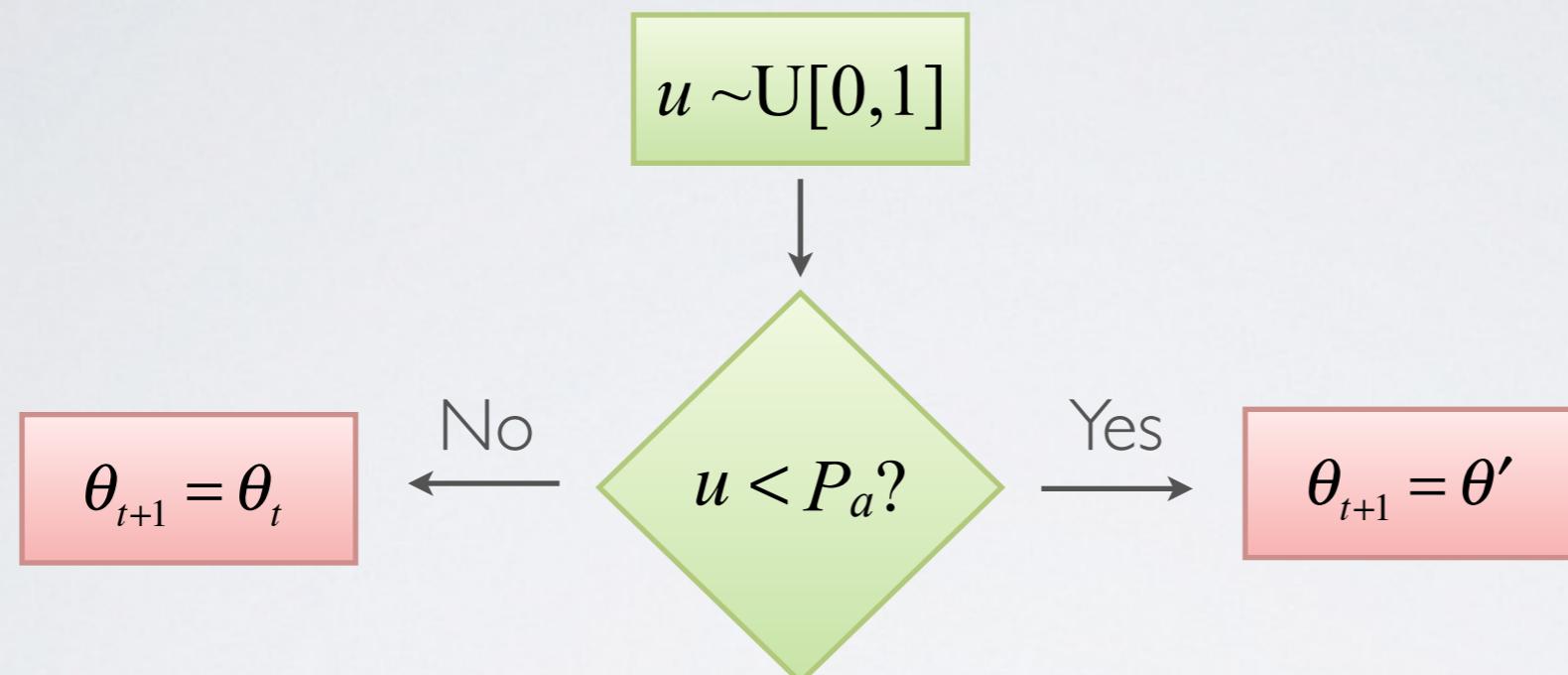
Accept/Reject as a Hypothesis Testing

$$P_a = \min \left\{ 1, \frac{q(\theta_t | \theta') S_0(\theta')}{q(\theta' | \theta_t) S_0(\theta_t)} \right\}$$



Accept/Reject as a Hypothesis Testing

$$P_a = \min \left\{ 1, \frac{q(\theta_t | \theta') S_0(\theta')}{q(\theta' | \theta_t) S_0(\theta_t)} \right\}$$



Accept/Reject as a Hypothesis Testing

$$u < \min \left\{ 1, \frac{q(\theta_t | \theta') S_0(\theta')}{q(\theta' | \theta_t) S_0(\theta_t)} \right\}$$

Accept/Reject as a Hypothesis Testing

$$u < \frac{q(\theta_t | \theta') S_0(\theta')}{q(\theta' | \theta_t) S_0(\theta_t)}$$

Accept/Reject as a Hypothesis Testing

$$u < \frac{q(\theta_t | \theta') S_0(\theta')}{q(\theta' | \theta_t) S_0(\theta_t)}$$

$$u < \frac{q(\theta_t | \theta') P(\theta')}{q(\theta' | \theta_t) P(\theta_t)} \cdot \prod_{i=1}^N \frac{P(x_i | \theta')}{P(x_i | \theta_t)}$$

Accept/Reject as a Hypothesis Testing

$$u < \frac{q(\theta_t | \theta') S_0(\theta')}{q(\theta' | \theta_t) S_0(\theta_t)}$$

$$u < \frac{q(\theta_t | \theta') P(\theta')}{q(\theta' | \theta_t) P(\theta_t)} \cdot \prod_{i=1}^N \frac{P(x_i | \theta')}{P(x_i | \theta_t)}$$

$$\log(u) < \log\left(\frac{q(\theta_t | \theta') P(\theta')}{q(\theta' | \theta_t) P(\theta_t)}\right) + \sum_{i=1}^N \log\left(\frac{P(x_i | \theta')}{P(x_i | \theta_t)}\right)$$

Accept/Reject as a Hypothesis Testing

$$u < \frac{q(\theta_t | \theta') S_0(\theta')}{q(\theta' | \theta_t) S_0(\theta_t)}$$

$$u < \frac{q(\theta_t | \theta') P(\theta')}{q(\theta' | \theta_t) P(\theta_t)} \cdot \prod_{i=1}^N \frac{P(x_i | \theta')}{P(x_i | \theta_t)}$$

$$\log(u) < \log\left(\frac{q(\theta_t | \theta') P(\theta')}{q(\theta' | \theta_t) P(\theta_t)}\right) + \sum_{i=1}^N \log\left(\frac{P(x_i | \theta')}{P(x_i | \theta_t)}\right)$$

$\frac{1}{N} \sum_{i=1}^N y_i > \mu_0 ?$

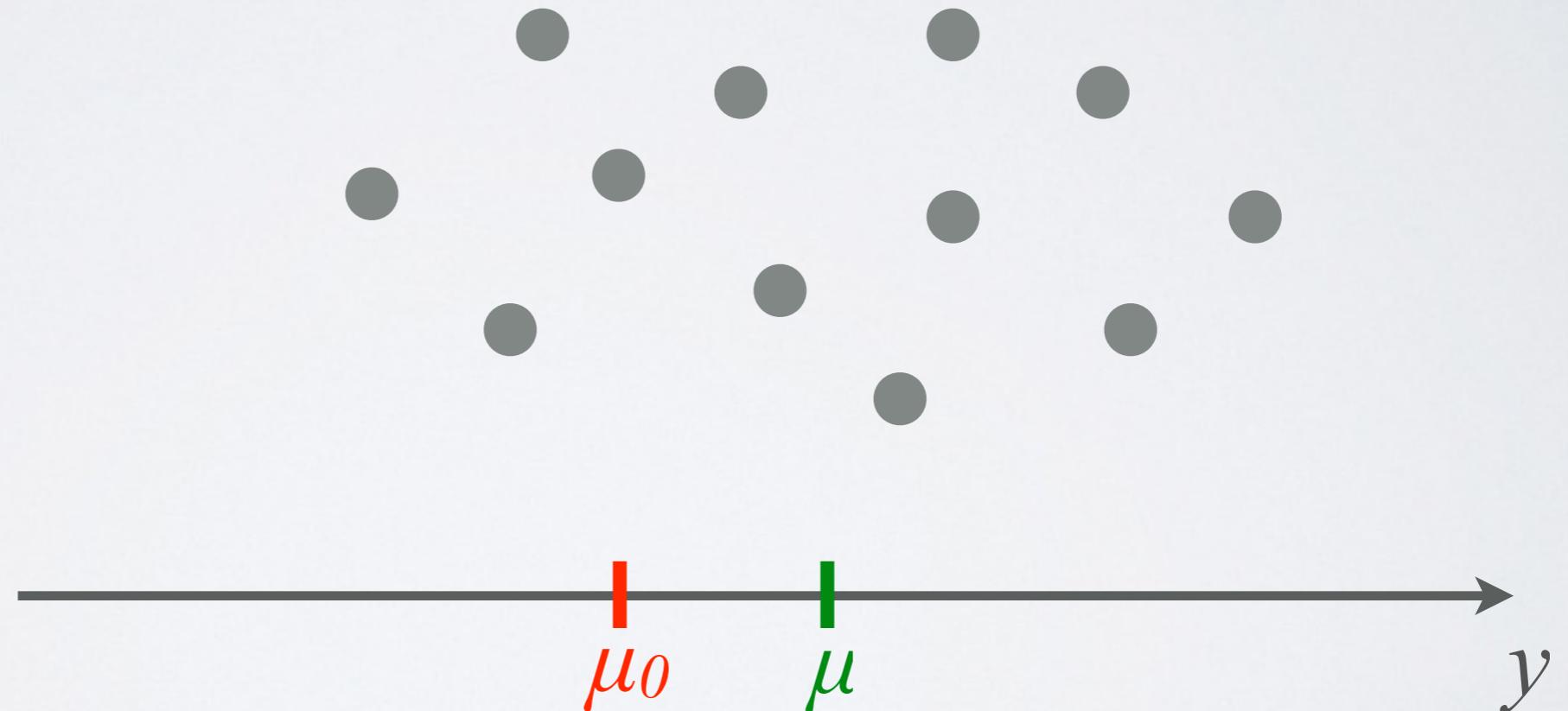
Accept/Reject as a Hypothesis Testing

- Given a set $\mathcal{Y} = \{y_i\}_{i=1}^N$, with a mean μ and a constant μ_0 , test:

$$\frac{1}{N} \sum_{i=1}^N y_i > \mu_0 ?$$

- $H_0: \mu > \mu_0$
- $H_a: \mu < \mu_0$

$$\log\left(\frac{P(x_i | \theta')}{P(x_i | \theta_t)}\right) \quad \frac{1}{N} \log\left(u \frac{q(\theta' | \theta_t) P(\theta_t)}{q(\theta_t | \theta') P(\theta')}\right)$$



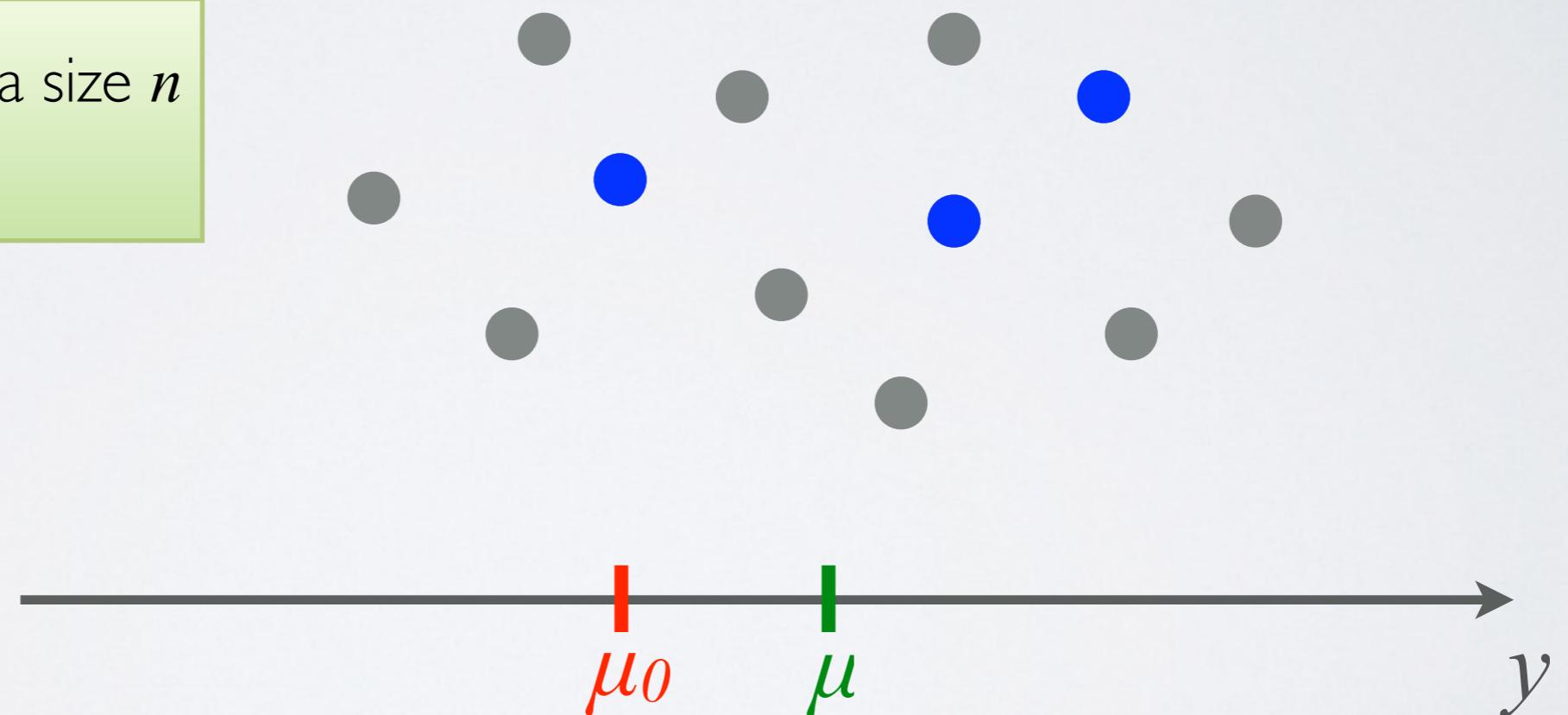
Accept/Reject as a Hypothesis Testing

- Given a set $\mathcal{Y} = \{y_i\}_{i=1}^N$, with a mean μ and a constant μ_0 , test:

$$\frac{1}{N} \sum_{i=1}^N y_i > \mu_0 ?$$

- $H_0: \mu > \mu_0$
- $H_a: \mu < \mu_0$

I. Draw $Y \subseteq \mathcal{Y}$ of a size n



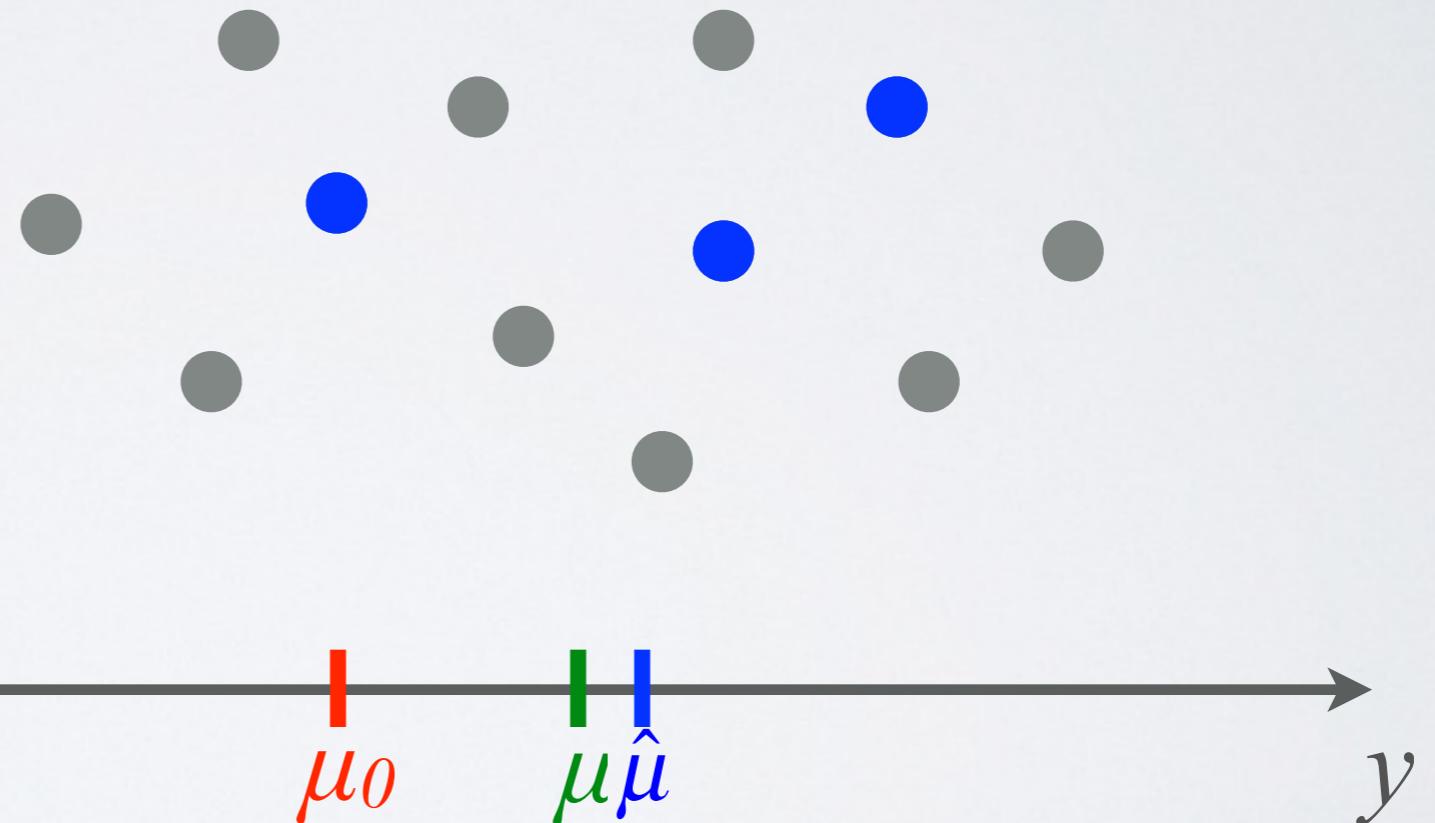
Accept/Reject as a Hypothesis Testing

- Given a set $\mathcal{Y} = \{y_i\}_{i=1}^N$, with a mean μ and a constant μ_0 , test:

$$\frac{1}{N} \sum_{i=1}^N y_i > \mu_0 ?$$

- $H_0: \mu > \mu_0$
- $H_a: \mu < \mu_0$

1. Draw $Y \subseteq \mathcal{Y}$ of a size n
2. Check $\hat{\mu} > \mu_0 ?$



Accept/Reject as a Hypothesis Testing

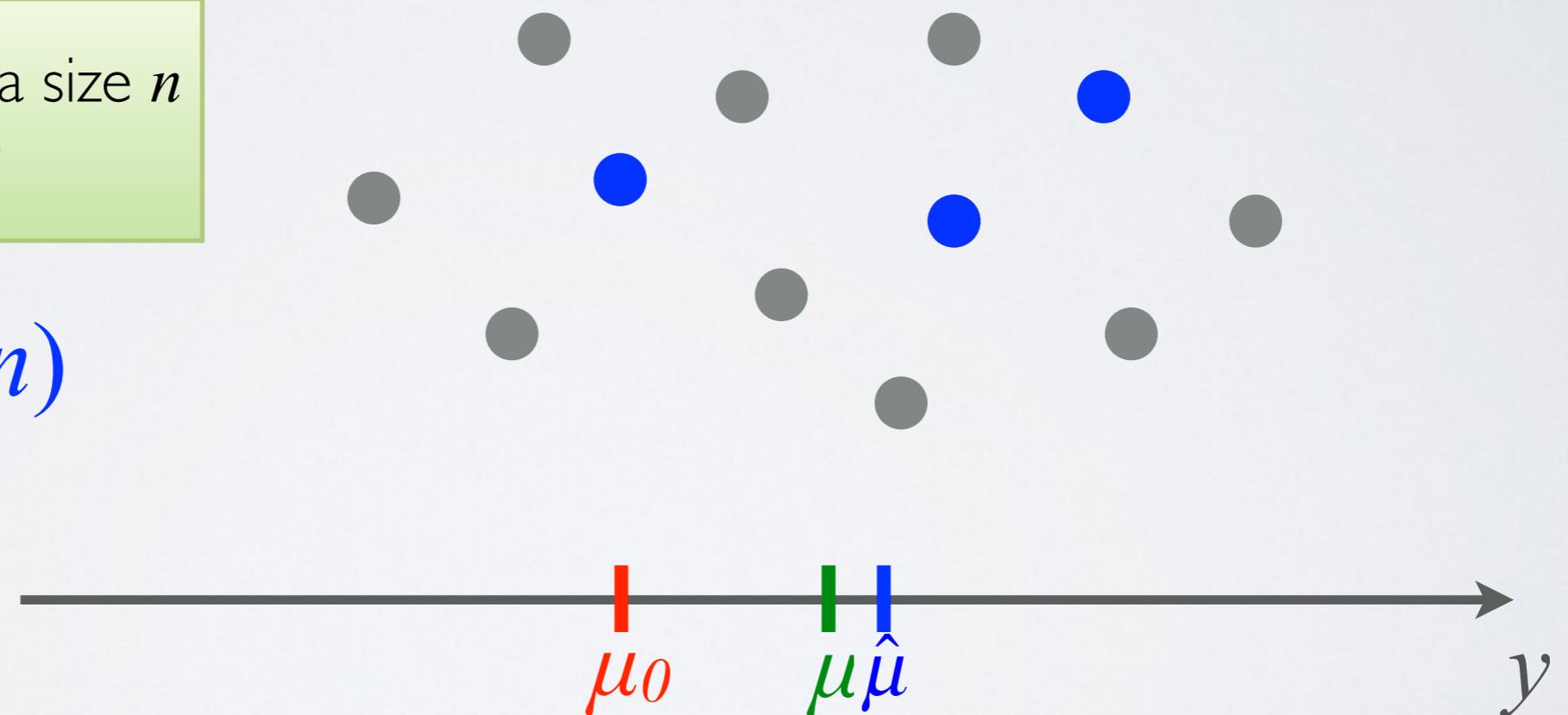
- Given a set $\mathcal{Y} = \{y_i\}_{i=1}^N$, with a mean μ and a constant μ_0 , test:

$$\frac{1}{N} \sum_{i=1}^N y_i > \mu_0 ?$$

- $H_0: \mu > \mu_0$
- $H_a: \mu < \mu_0$

1. Draw $Y \subseteq \mathcal{Y}$ of a size n
2. Check $\hat{\mu} > \mu_0 ?$

$O(N) \rightarrow O(n)$



Accept/Reject as a Hypothesis Testing

- Given a set $\mathcal{Y} = \{y_i\}_{i=1}^N$, with a mean μ and a constant μ_0 , test:

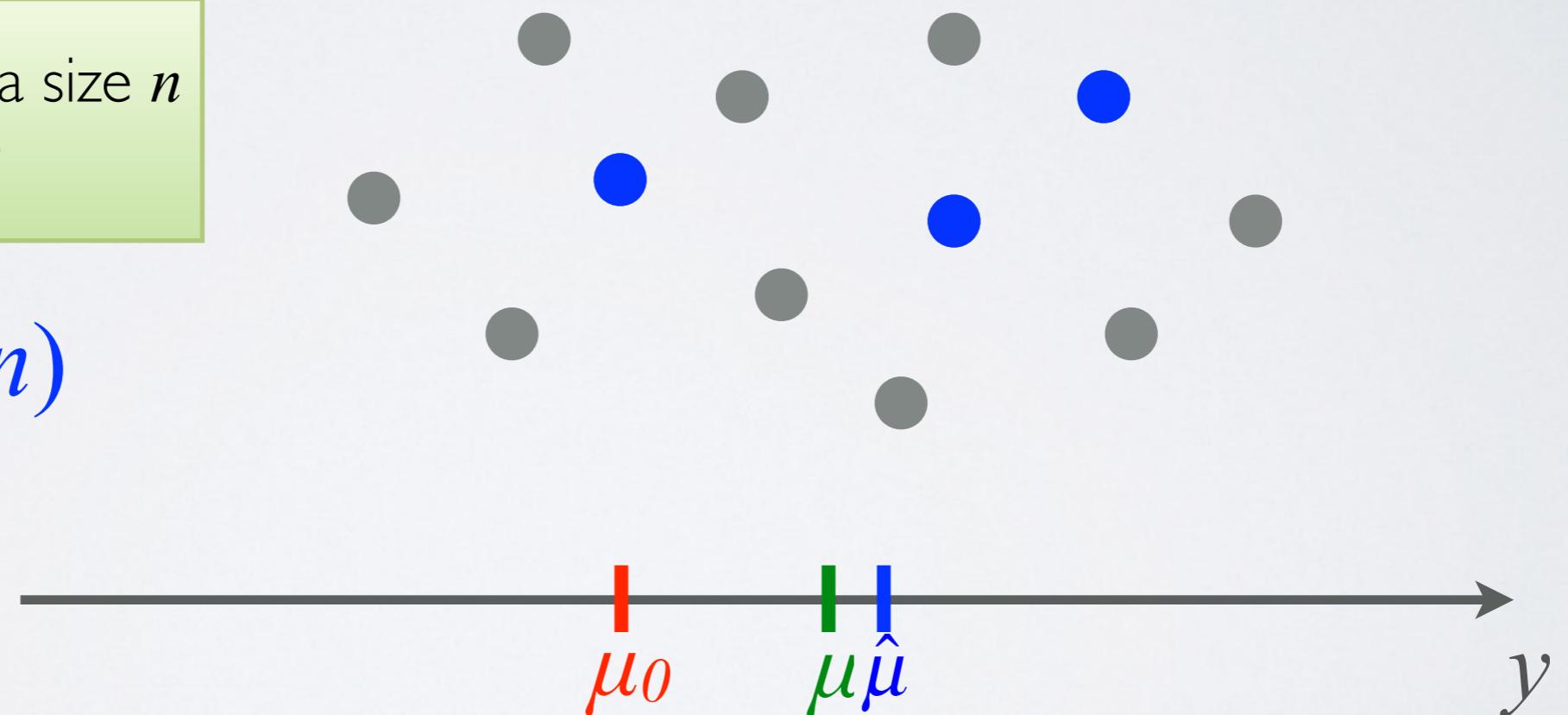
$$\frac{1}{N} \sum_{i=1}^N y_i > \mu_0 ?$$

- $H_0: \mu > \mu_0$
- $H_a: \mu < \mu_0$

1. Draw $Y \subseteq \mathcal{Y}$ of a size n
2. Check $\hat{\mu} > \mu_0 ?$

$O(N) \rightarrow O(n)$

$n?$



Accept/Reject as a Hypothesis Testing

- Given a set $\mathcal{Y} = \{y_i\}_{i=1}^N$, with a mean μ and a constant μ_0 , test:

$$\frac{1}{N} \sum_{i=1}^N y_i > \mu_0 ?$$

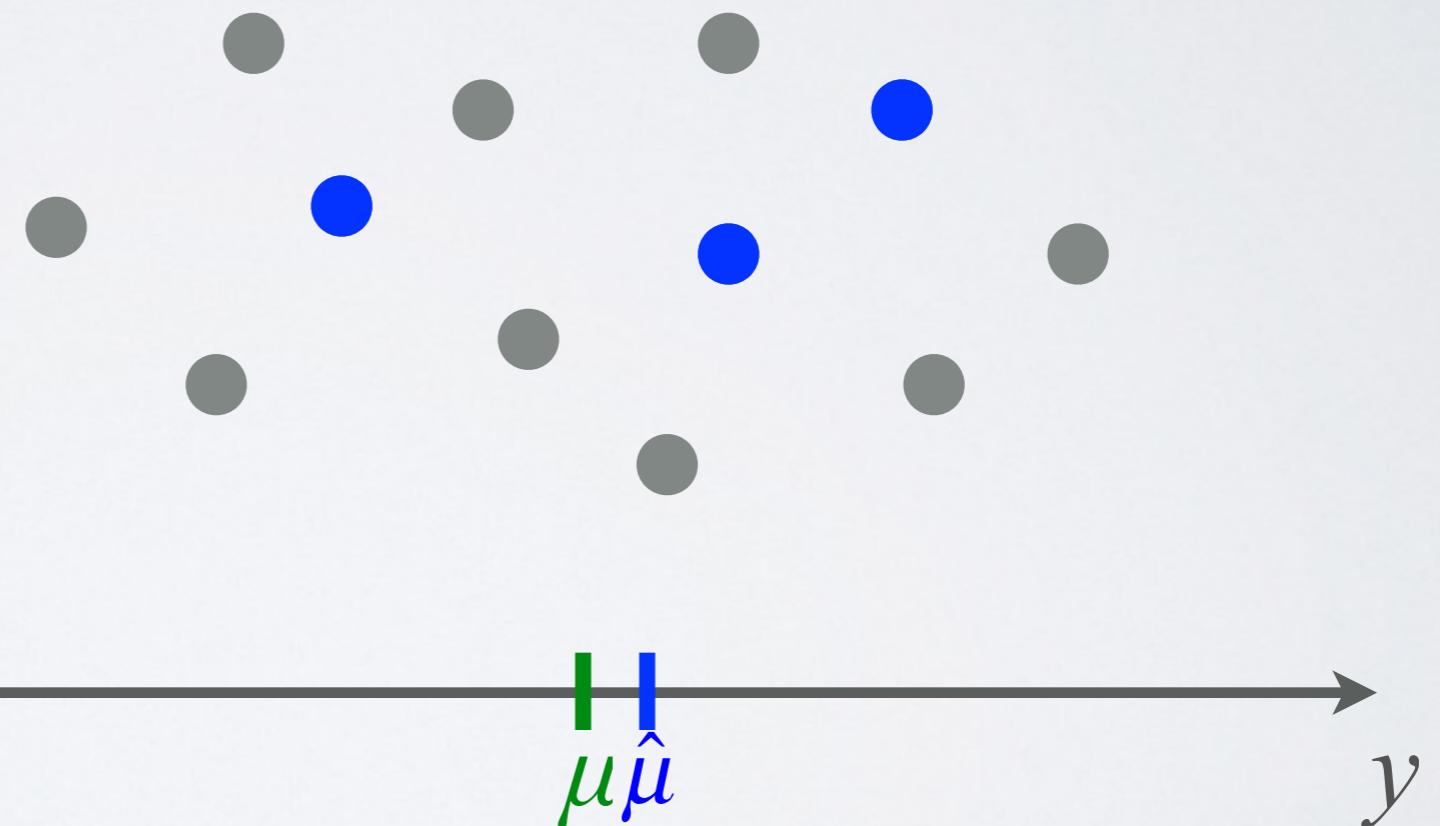
- $H_0: \mu > \mu_0$
- $H_a: \mu < \mu_0$

1. Draw $Y \subseteq \mathcal{Y}$ of a size n
2. Check $\hat{\mu} > \mu_0 ?$

$$\log\left(\frac{P(x_i | \theta')}{P(x_i | \theta_t)}\right) \quad \frac{1}{N} \log\left(u \frac{q(\theta' | \theta_t) P(\theta_t)}{q(\theta_t | \theta') P(\theta')}\right)$$

$O(N) \rightarrow O(n)$

$n?$
 μ_0
 $n \approx 1$



Accept/Reject as a Hypothesis Testing

- Given a set $\mathcal{Y} = \{y_i\}_{i=1}^N$, with a mean μ and a constant μ_0 , test:

$$\frac{1}{N} \sum_{i=1}^N y_i > \mu_0 ?$$

- $H_0: \mu > \mu_0$
- $H_a: \mu < \mu_0$

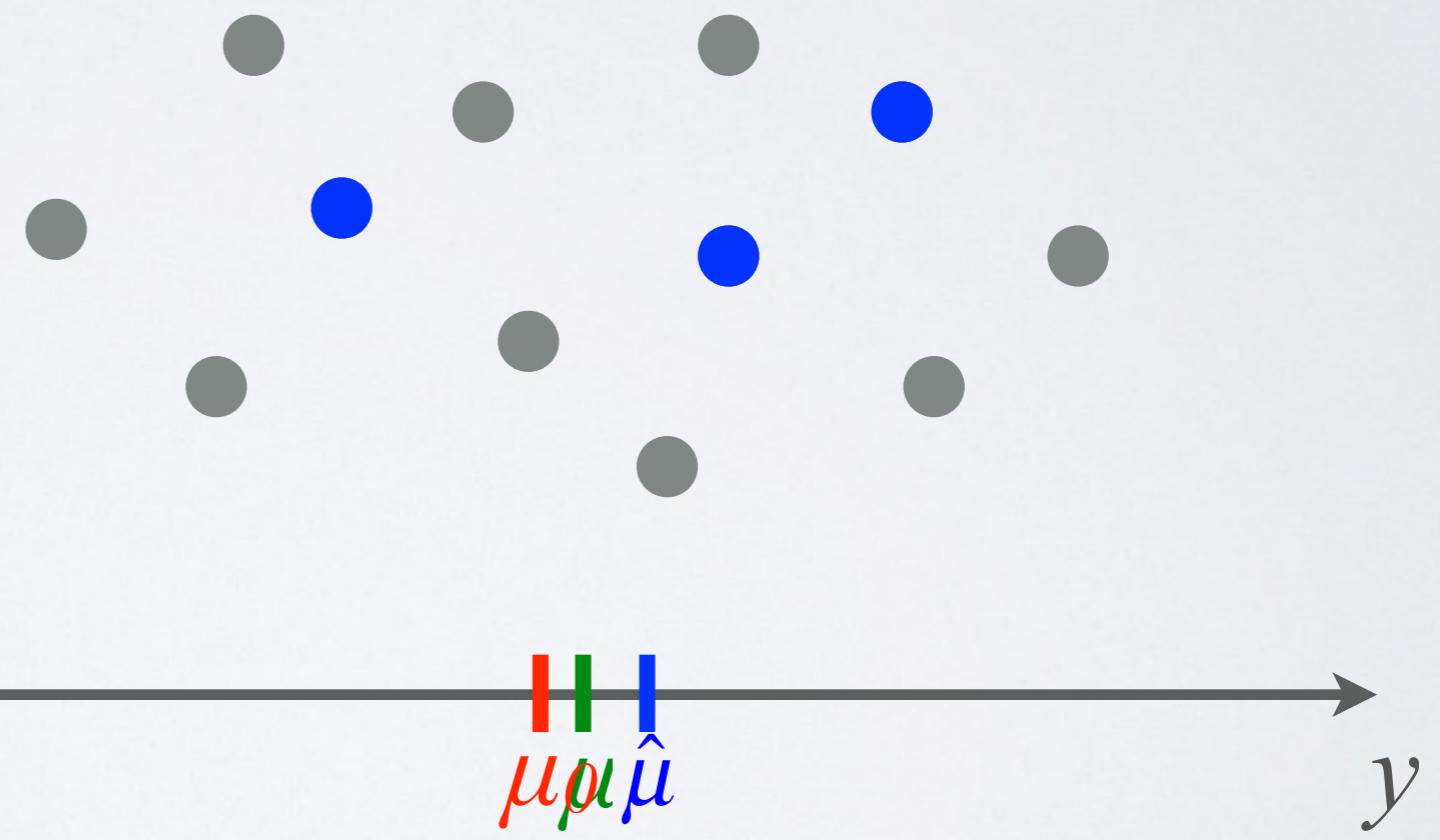
1. Draw $Y \subseteq \mathcal{Y}$ of a size n
2. Check $\hat{\mu} > \mu_0$?

$O(N) \rightarrow O(n)$

$n?$

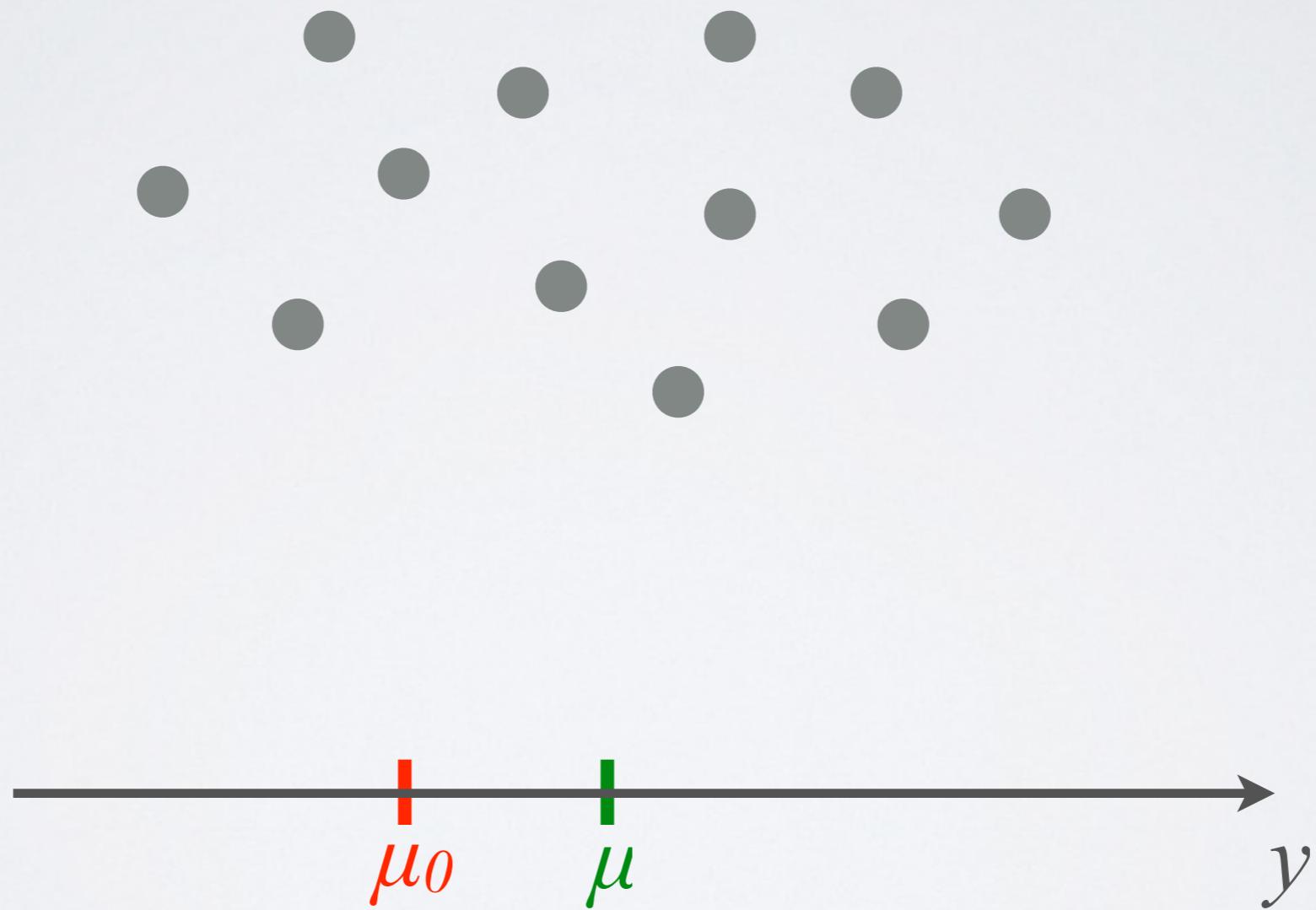
$n \approx 1$

$$\log\left(\frac{P(x_i | \theta')}{P(x_i | \theta_t)}\right) \quad \frac{1}{N} \log\left(u \frac{q(\theta' | \theta_t) P(\theta_t)}{q(\theta_t | \theta') P(\theta')}\right)$$



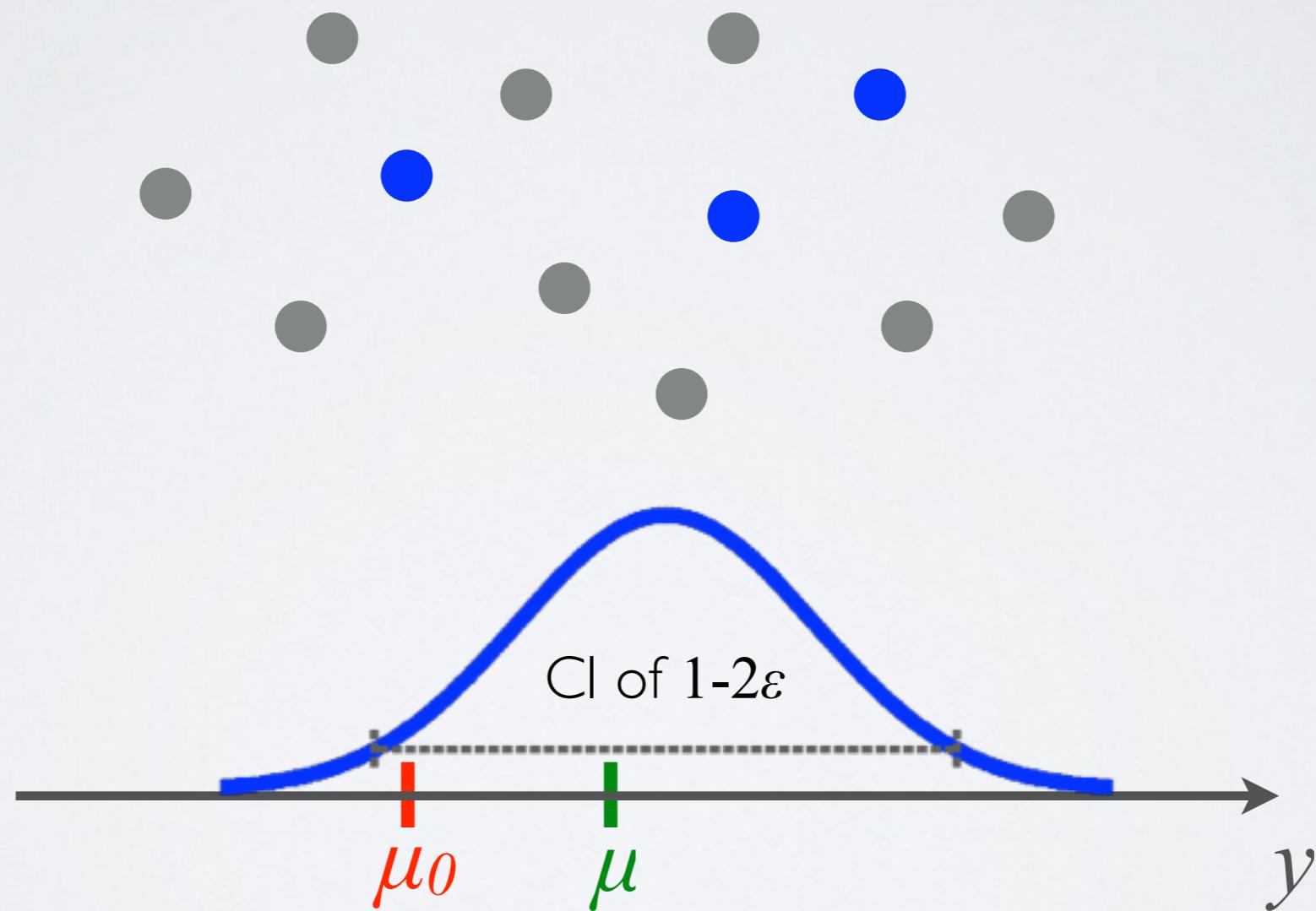
Sequential Hypothesis Testing

- Keep drawing data cases randomly *w/o replacement* until we are confident about our approximation.



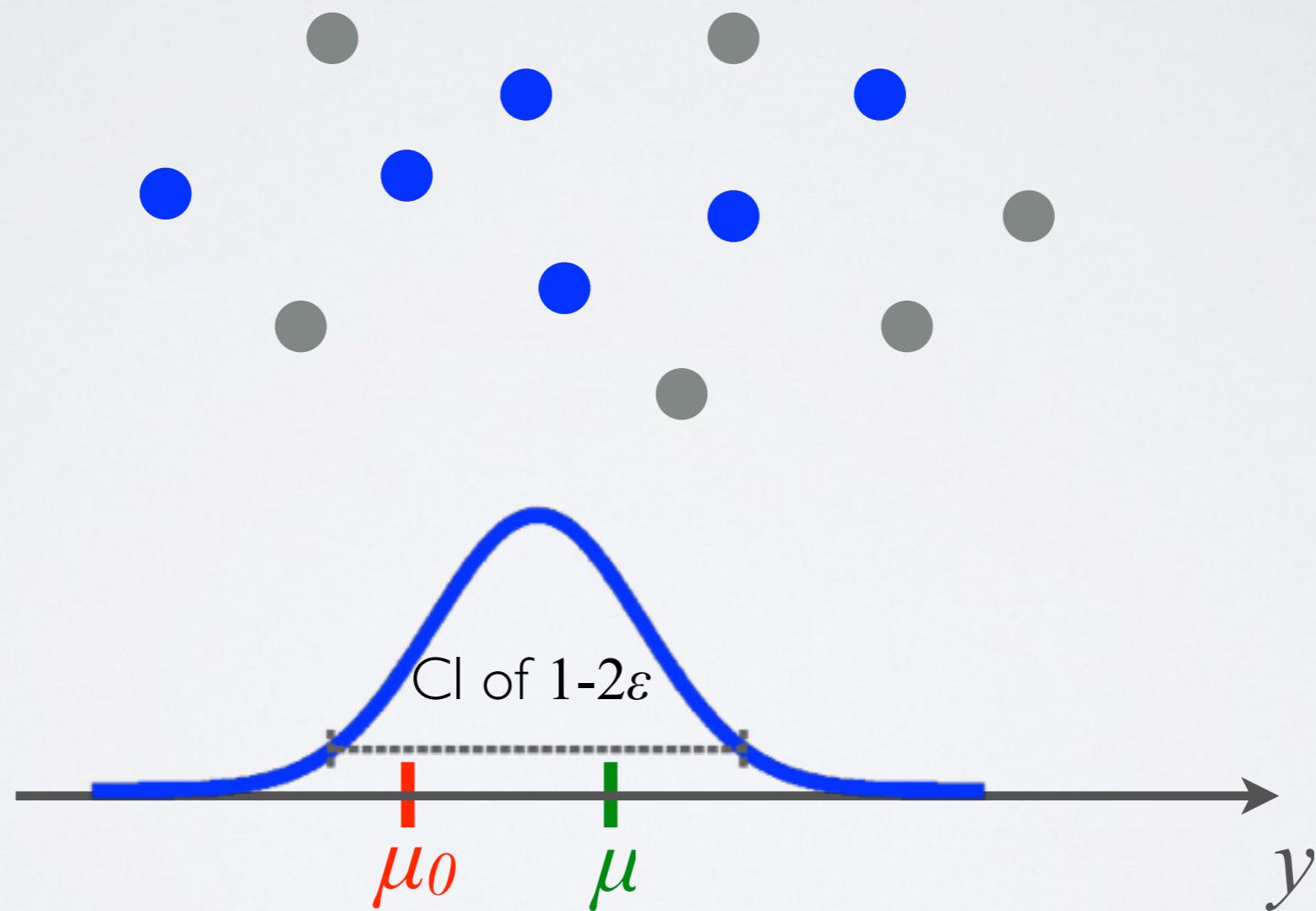
Sequential Hypothesis Testing

- Keep drawing data cases randomly *w/o replacement* until we are confident about our approximation.



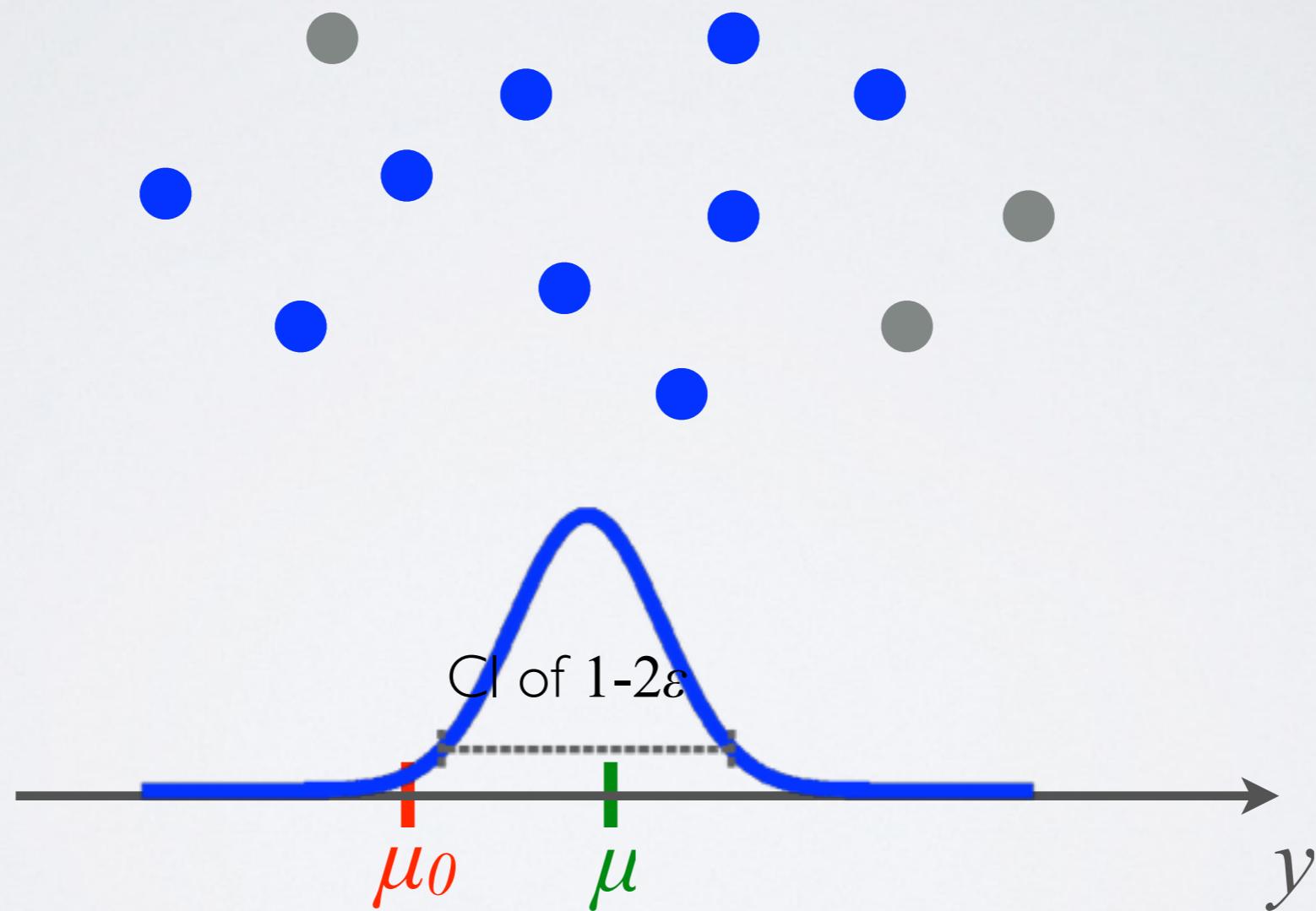
Sequential Hypothesis Testing

- Keep drawing data cases randomly *w/o replacement* until we are confident about our approximation.



Sequential Hypothesis Testing

- Keep drawing data cases randomly *w/o replacement* until we are confident about our approximation.



Sequential Hypothesis Testing

Sequential Hypothesis Testing

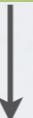
Draw $Y_j \subseteq \mathcal{Y}$ w/o replacement

$$n \leftarrow n + m$$

Sequential Hypothesis Testing

Draw $Y_j \subseteq \mathcal{Y}$ w/o replacement

$$n \leftarrow n+m$$



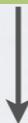
$$\hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n y_i, \quad \hat{\sigma}_j = \frac{\hat{\sigma}(y_i)}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}$$

Finite population
correction factor

Sequential Hypothesis Testing

Draw $Y_j \subseteq \mathcal{Y}$ w/o replacement

$$n \leftarrow n+m$$

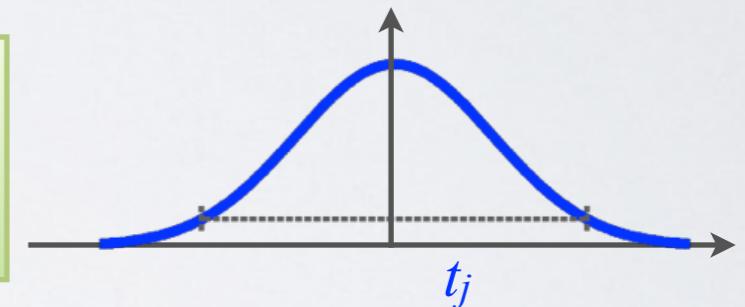


$$\hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n y_i, \quad \hat{\sigma}_j = \frac{\hat{\sigma}(y_i)}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}$$

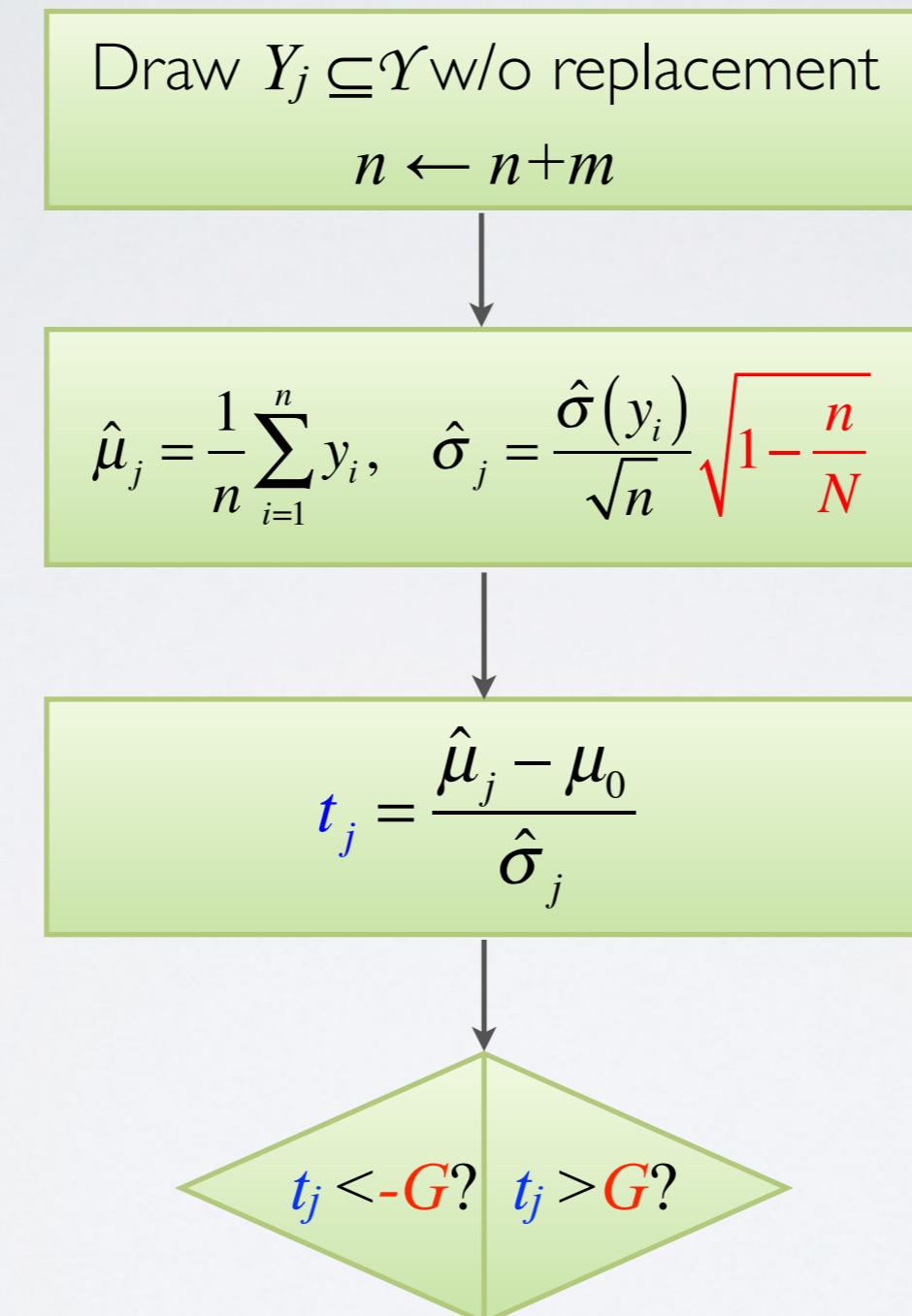


$$t_j = \frac{\hat{\mu}_j - \mu_0}{\hat{\sigma}_j}$$

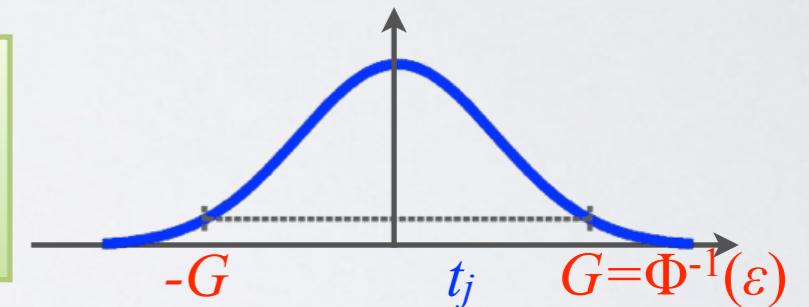
Finite population
correction factor



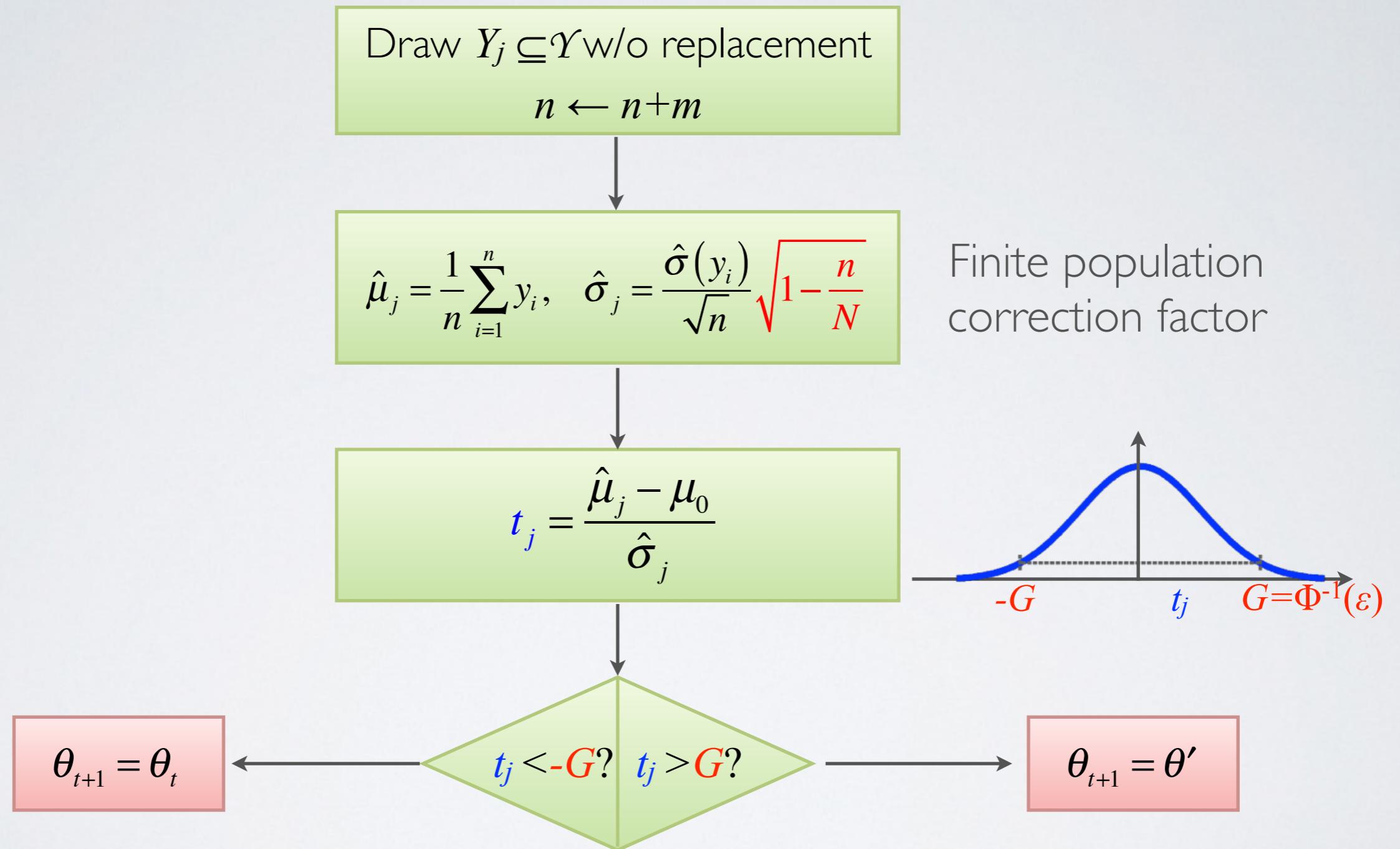
Sequential Hypothesis Testing



Finite population
correction factor



Sequential Hypothesis Testing



Sequential Hypothesis Testing

Need more data

$$j \leftarrow j+1$$

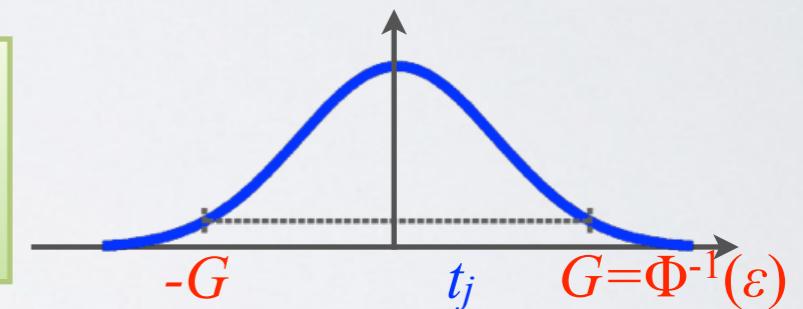
Draw $Y_j \subseteq \mathcal{Y}$ w/o replacement

$$n \leftarrow n+m$$

$$\hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n y_i, \quad \hat{\sigma}_j = \frac{\hat{\sigma}(y_i)}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}$$

Finite population
correction factor

$$t_j = \frac{\hat{\mu}_j - \mu_0}{\hat{\sigma}_j}$$

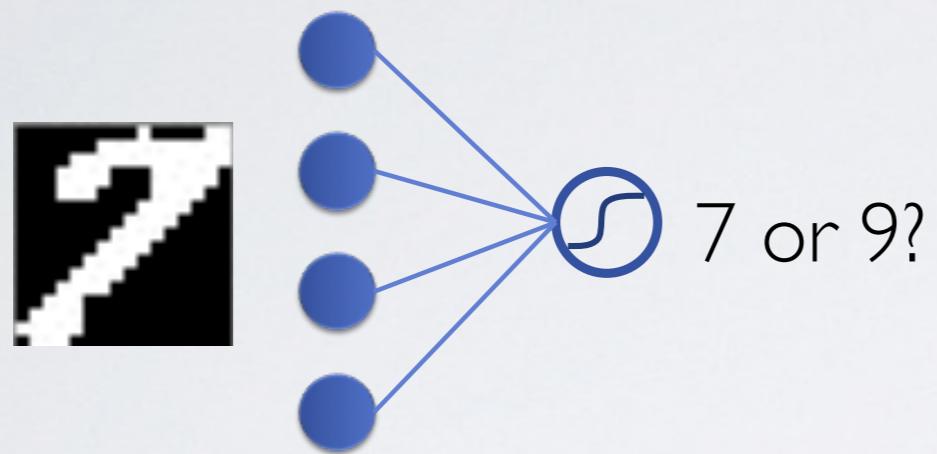


$$\theta_{t+1} = \theta_t$$

$$t_j < -G? \quad t_j > G?$$

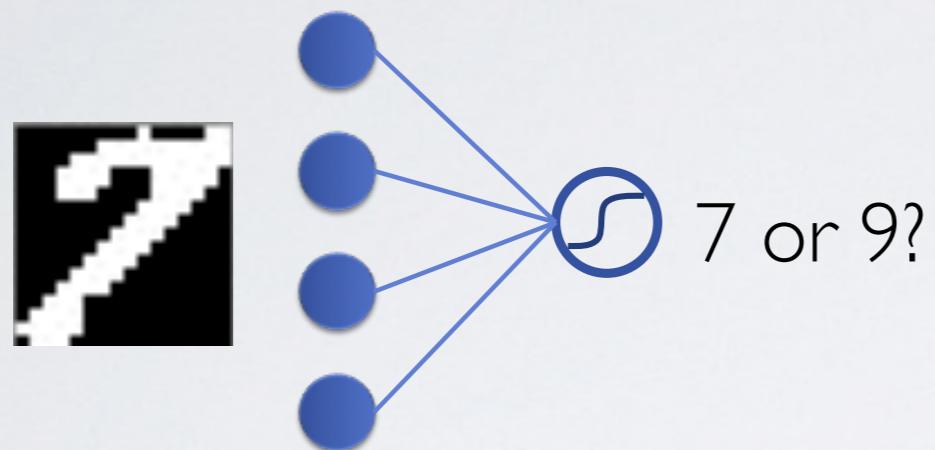
$$\theta_{t+1} = \theta'$$

Application I: Posterior Inference for Bayesian Logistic Regression



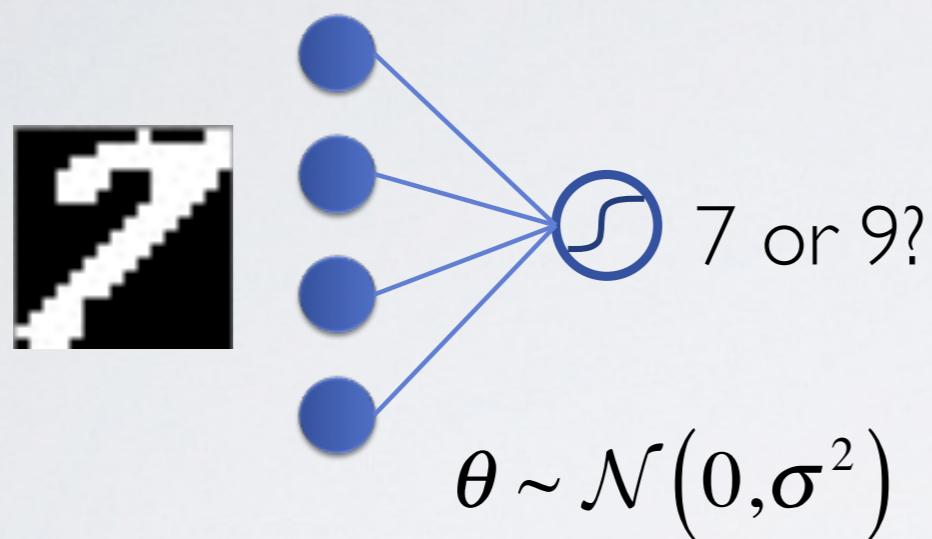
Application I: Posterior Inference for Bayesian Logistic Regression

- Data: MNIST digit 7 and 9, 12214 data points.



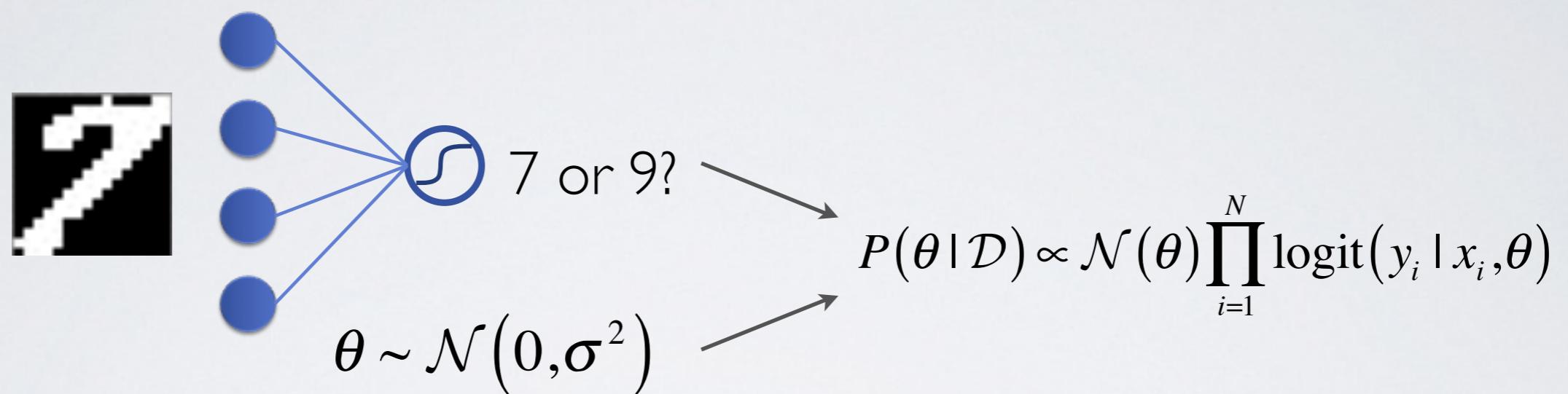
Application I: Posterior Inference for Bayesian Logistic Regression

- Data: MNIST digit 7 and 9, 12214 data points.



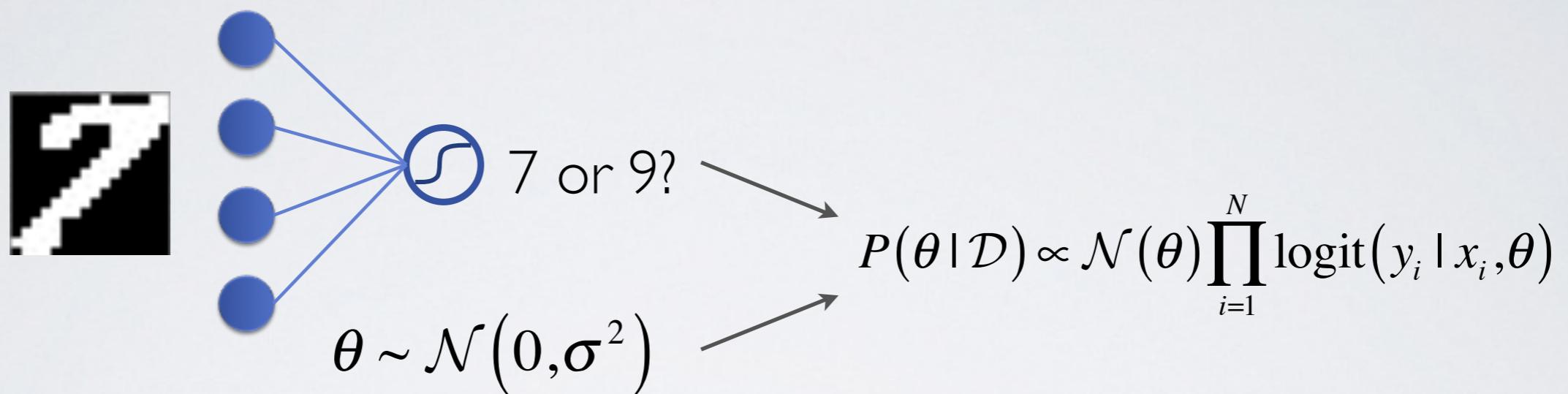
Application I: Posterior Inference for Bayesian Logistic Regression

- Data: MNIST digit 7 and 9, 12214 data points.



Application I: Posterior Inference for Bayesian Logistic Regression

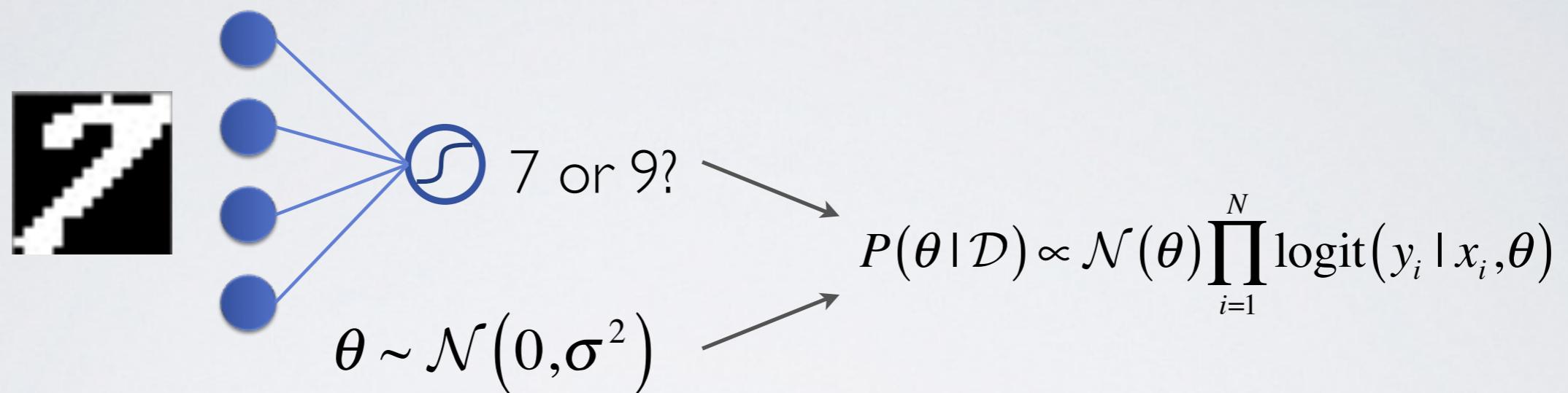
- Data: MNIST digit 7 and 9, 12214 data points.



- Sampler: M-H with Random walk proposal.

Application I: Posterior Inference for Bayesian Logistic Regression

- Data: MNIST digit 7 and 9, 12214 data points.

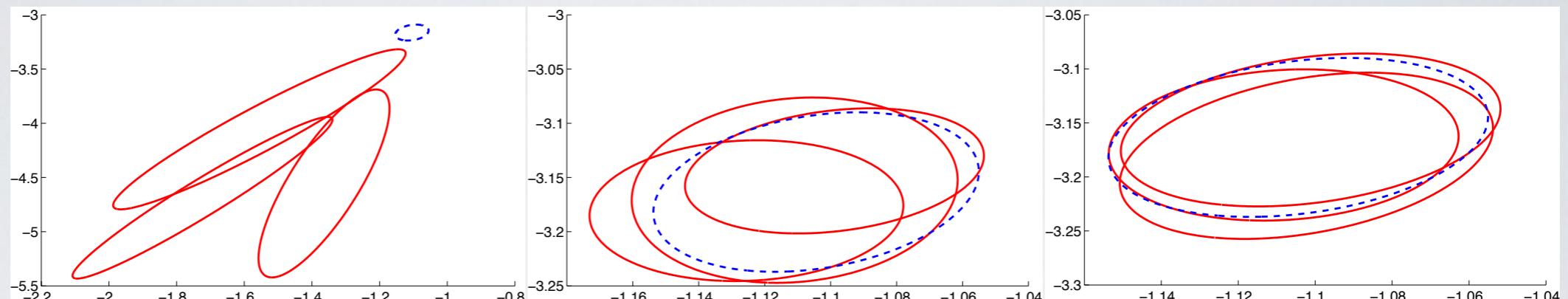


- Sampler: M-H with Random walk proposal.
- Quantity of interest: $P(x^* | \mathcal{D}) = \mathbb{E}_{P(\theta|\mathcal{D})} [\text{logit}(x^* | \theta)]$

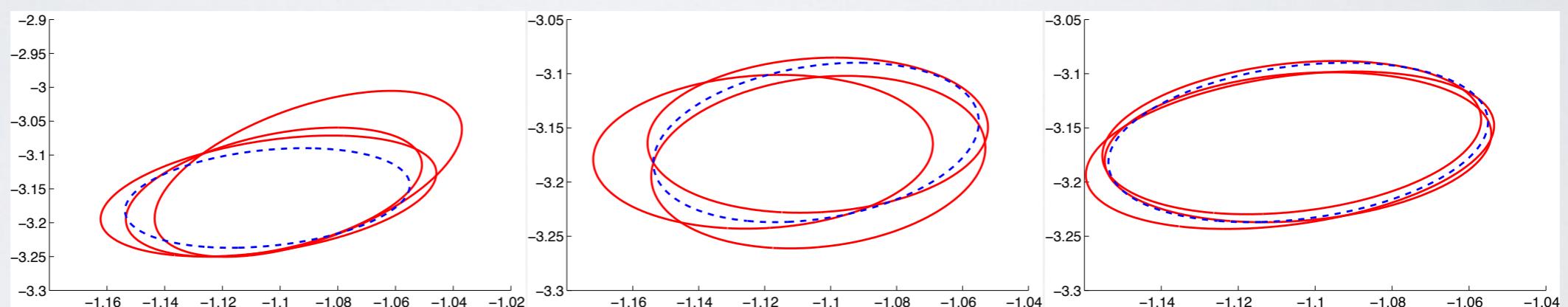
$$\downarrow$$

$$f(\theta)$$

Exact MH



Approximate
 $\varepsilon=0.01$



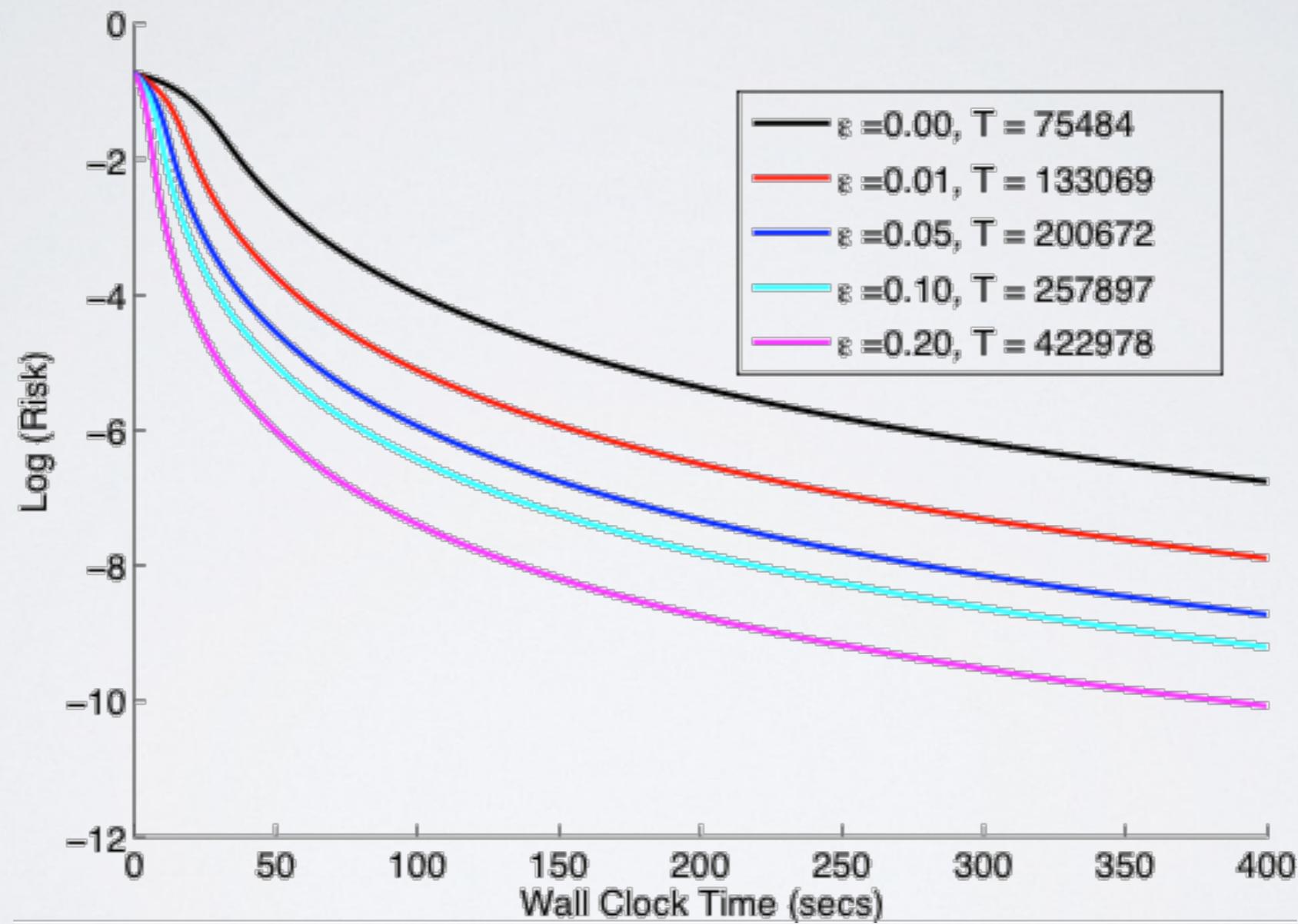
50 sec

100 sec

400 sec

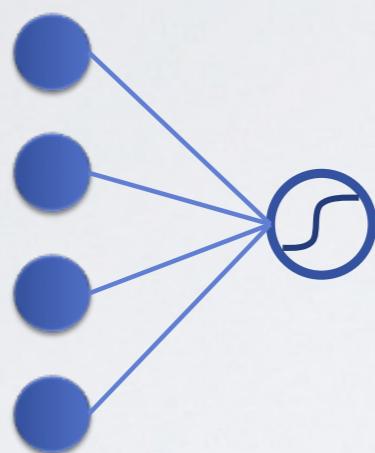
Wall Clock Time

$$Risk = \mathbb{E} \left[\left(\langle f \rangle_{S_0} - \overline{f(\theta_t)} \right)^2 \right]$$



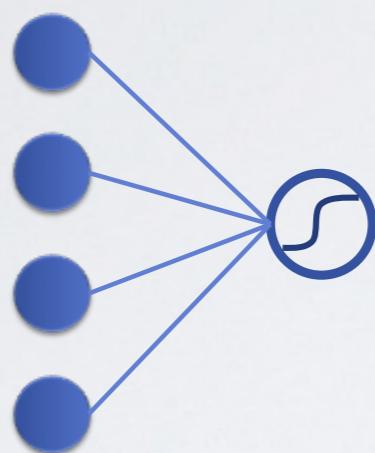
Application II: Variable selection in Logistic Regression

- Data: MiniBooNE | 30,065 data points with 50 features
- Sampler: Reversible jump MCMC



Application II: Variable selection in Logistic Regression

- Data: MiniBooNE | 30,065 data points with 50 features
- Sampler: Reversible jump MCMC



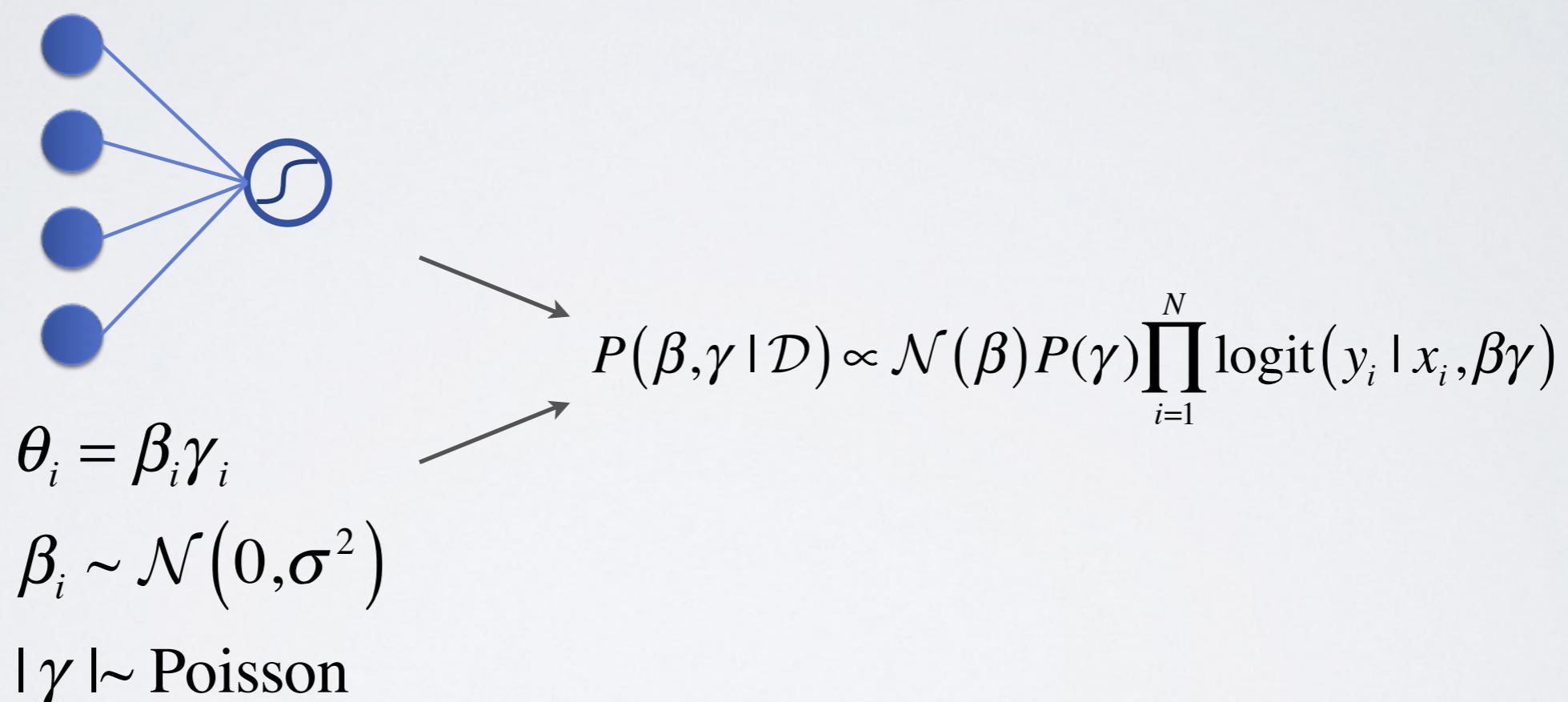
$$\theta_i = \beta_i \gamma_i$$

$$\beta_i \sim \mathcal{N}(0, \sigma^2)$$

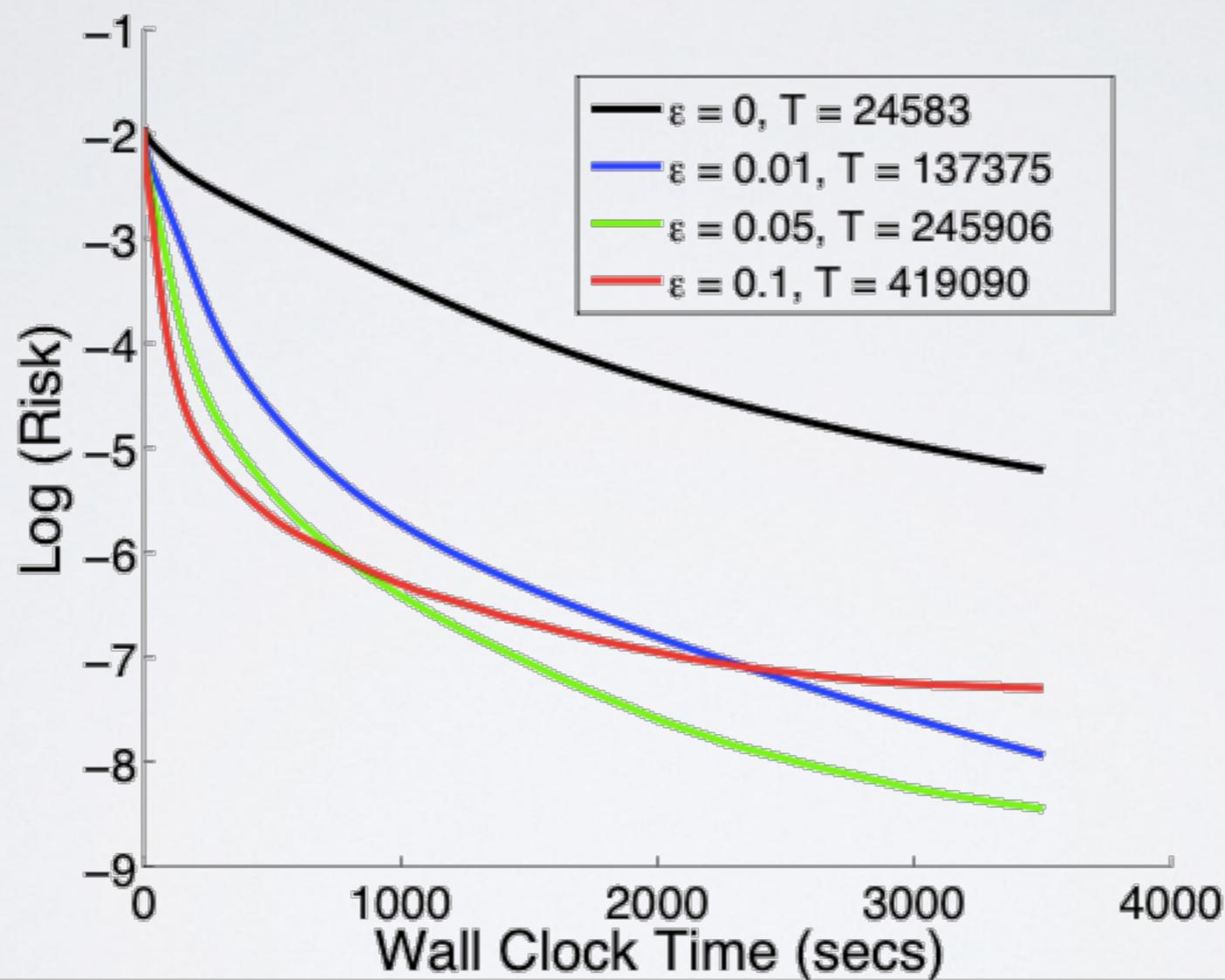
$$|\gamma| \sim \text{Poisson}$$

Application II: Variable selection in Logistic Regression

- Data: MiniBooNE | 30,065 data points with 50 features
- Sampler: Reversible jump MCMC

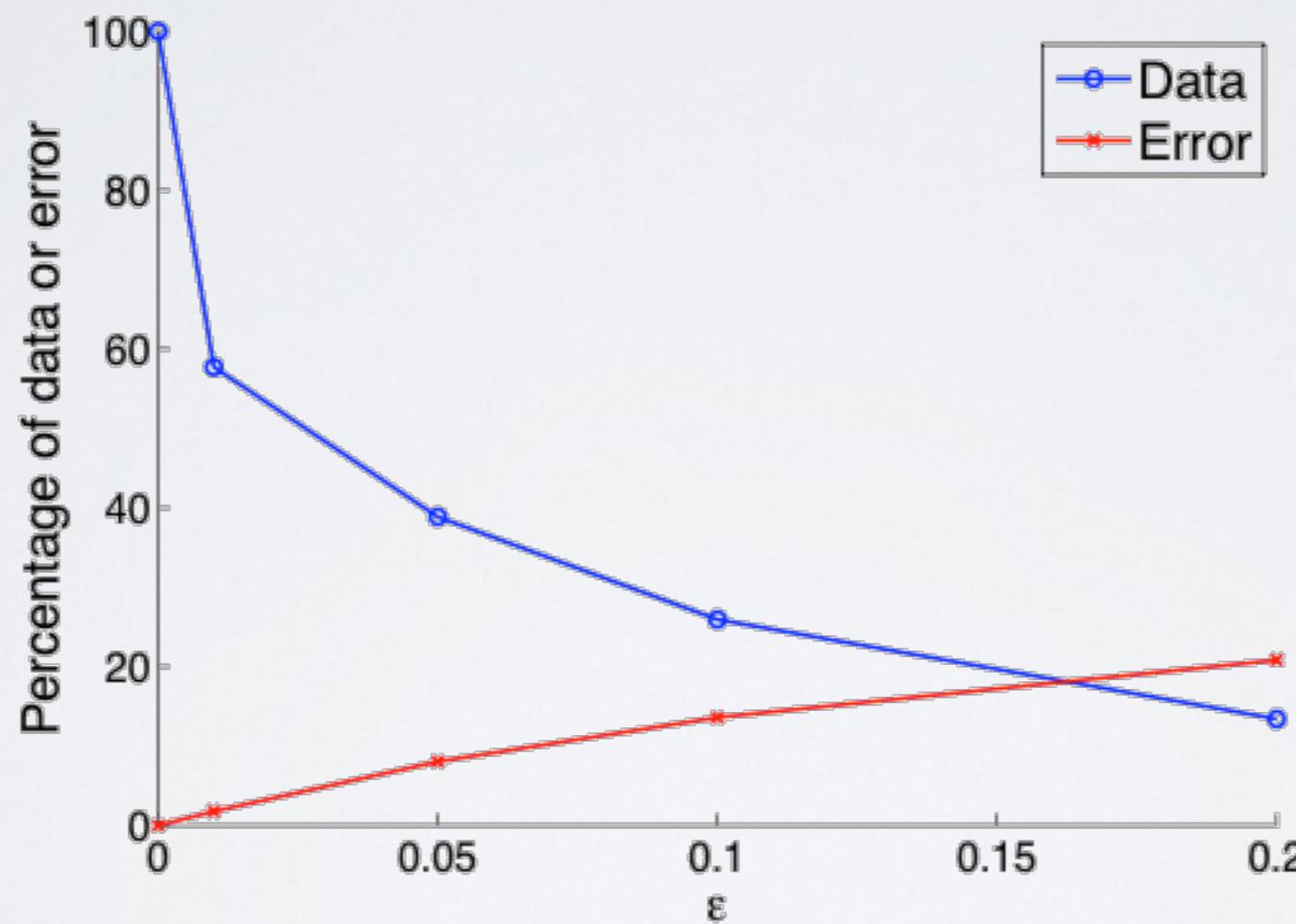


$$Risk = \mathbb{E} \left[\left(\langle f \rangle_{S_0} - \overline{f(\theta_t)} \right)^2 \right]$$

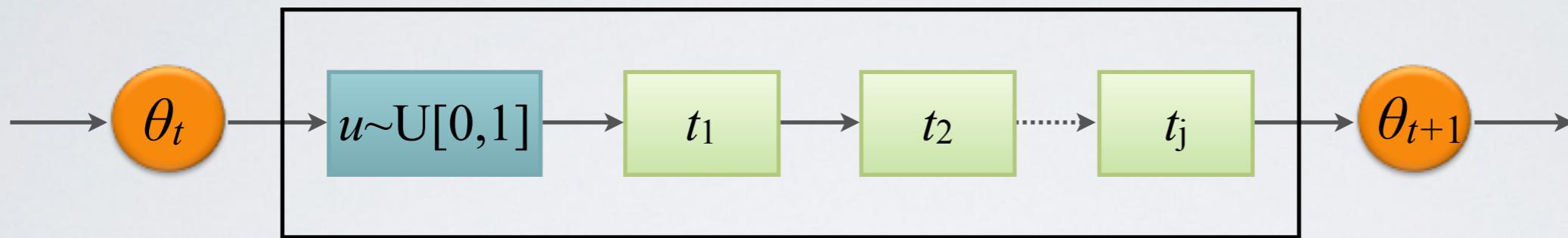


How to Design the Sequential Test?

$\mu > \mu_0?$

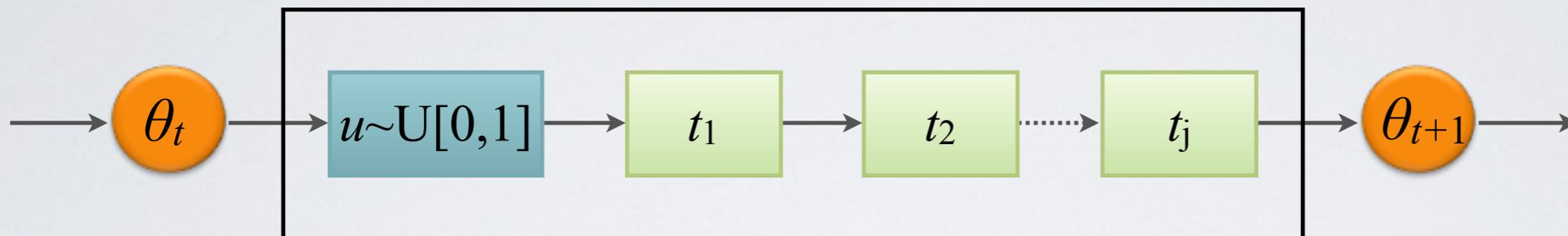


Sequential Test Design



Sequential Test Design

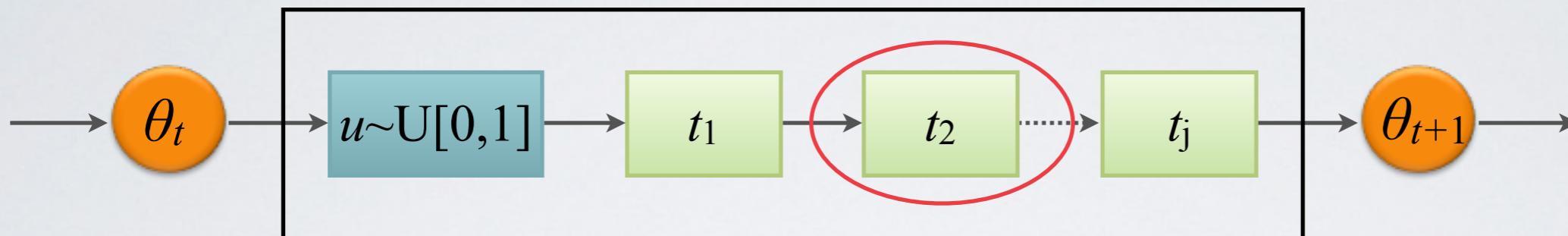
- How to design the sequential test to minimize the *usage of data* subject to a tolerance of the **error**?



- The stationary distribution, S_ε

Sequential Test Design

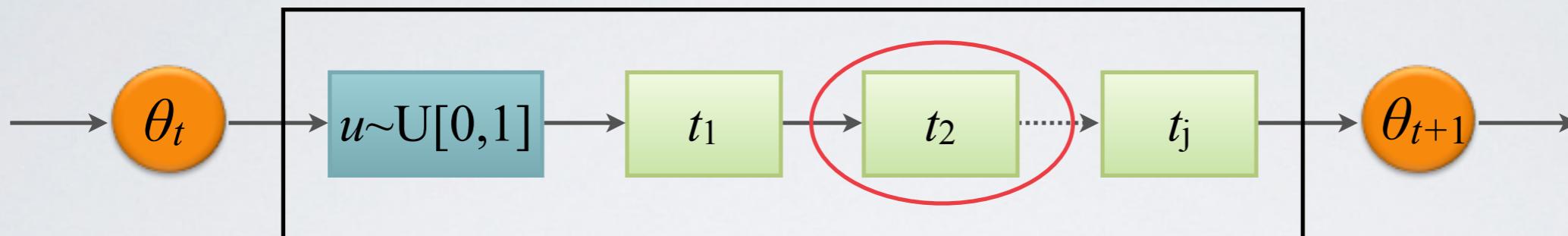
- How to design the sequential test to minimize the *usage of data* subject to a tolerance of the **error**?



- One step of test, error: ε , data: m
- The stationary distribution, S_ε

Sequential Test Design

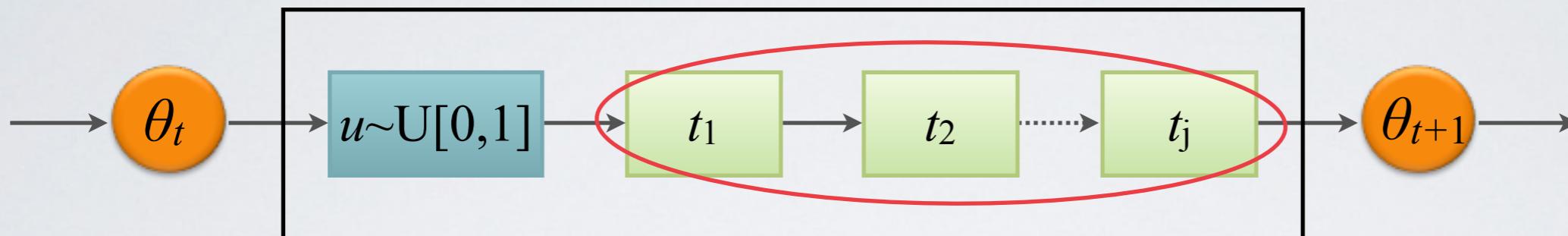
- How to design the sequential test to minimize the *usage of data* subject to a tolerance of the **error**?



- One step of test, error: ε , data: m
- The stationary distribution, S_ε

Sequential Test Design

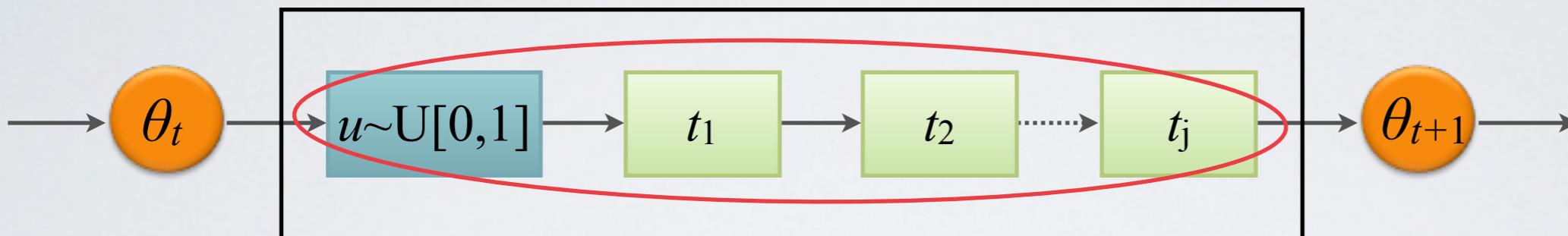
- How to design the sequential test to minimize the *usage of data* subject to a tolerance of the **error**?



- One step of test, error: ε , data: m
- One sequential test, conditioned on u
- The stationary distribution, S_ε

Sequential Test Design

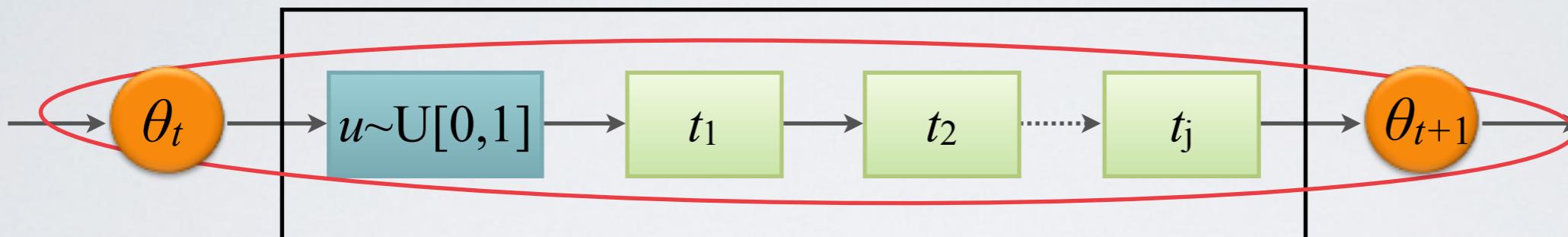
- How to design the sequential test to minimize the *usage of data* subject to a tolerance of the **error**?



- One step of test, error: ε , data: m
- One sequential test, conditioned on u
- One accept/reject step (marginalized over u)
- The stationary distribution, S_ε

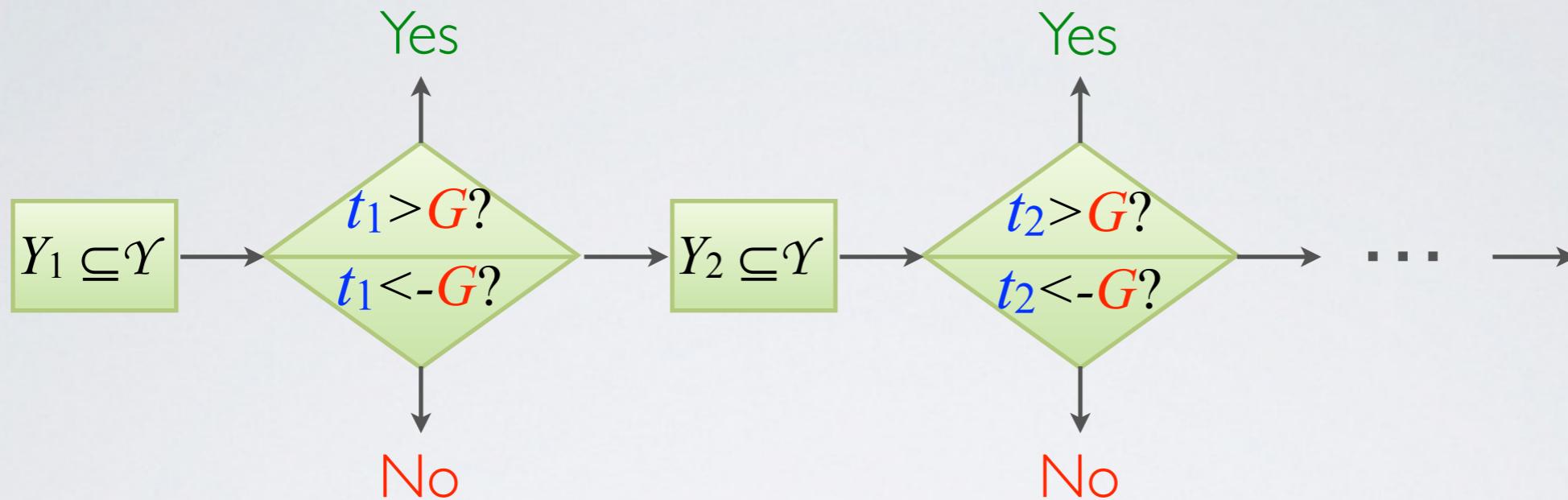
Sequential Test Design

- How to design the sequential test to minimize the *usage of data* subject to a tolerance of the **error**?



- One step of test, error: ε , data: m
- One sequential test, conditioned on u
- One accept/reject step (marginalized over u)
- The stationary distribution, S_ε

One Sequential Test



$\mu > \mu_0$

$$\bar{\pi} = \mathbb{E}[j^*] / J$$

$$\mathcal{E} = P(t_{j^*} < -G)$$

Random Walk of t_j

- Two assumptions about Y_j , when m is large $Y_1 \subseteq \gamma \rightarrow Y_2 \subseteq \gamma \rightarrow \dots$
 1. $\{\bar{y}_j\}$ is jointly Normal
 2. Accurate estimate of σ_y , $t_j = \frac{\overline{\{\bar{y}_j\}} - \mu_0}{\hat{\sigma}_j} \approx \frac{\overline{\{\bar{y}_j\}} - \mu_0}{\sigma_j} \sim \mathcal{N}(t_j)$

Random Walk of t_j

- Two assumptions about Y_j , when m is large $Y_1 \subseteq \gamma \rightarrow Y_2 \subseteq \gamma \rightarrow \dots$
 1. $\{\bar{y}_j\}$ is jointly Normal
 2. Accurate estimate of σ_y , $t_j = \frac{\overline{\{\bar{y}_j\}} - \mu_0}{\hat{\sigma}_j} \approx \frac{\overline{\{\bar{y}_j\}} - \mu_0}{\sigma_j} \sim \mathcal{N}(t_j)$

$$P(\bar{y}_1, \bar{y}_2, \dots, \bar{y}_j) = \mathcal{N}(\mu, \Sigma)$$

Random Walk of t_j

- Two assumptions about Y_j , when m is large $Y_1 \subseteq \gamma \rightarrow Y_2 \subseteq \gamma \rightarrow \dots$
 1. $\{\bar{y}_j\}$ is jointly Normal
 2. Accurate estimate of σ_y , $t_j = \frac{\overline{\{\bar{y}_j\}} - \mu_0}{\hat{\sigma}_j} \approx \frac{\overline{\{\bar{y}_j\}} - \mu_0}{\sigma_j} \sim \mathcal{N}(t_j)$

$$P(\bar{y}_1, \bar{y}_2, \dots, \bar{y}_j) = \mathcal{N}(\mu, \Sigma)$$

- Sample with Replacement

$$\Sigma(\bar{y}) = \frac{\sigma_y^2}{m} \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix}$$

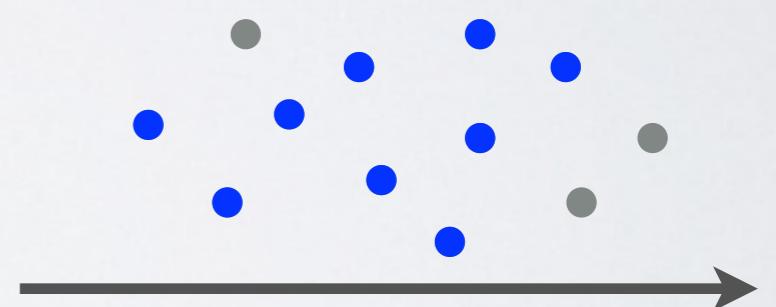
Random Walk of t_j

- Two assumptions about Y_j , when m is large $Y_1 \subseteq \gamma \rightarrow Y_2 \subseteq \gamma \rightarrow \dots$
 - $\{\bar{y}_j\}$ is jointly Normal
 - Accurate estimate of σ_y , $t_j = \frac{\{\bar{y}_j\} - \mu_0}{\hat{\sigma}_j} \approx \frac{\{\bar{y}_j\} - \mu_0}{\sigma_j} \sim \mathcal{N}(t_j)$

$$P(\bar{y}_1, \bar{y}_2, \dots, \bar{y}_j) = \mathcal{N}(\mu, \Sigma)$$

- Sample without Replacement

$$\Sigma(\bar{y}) = \frac{\sigma^2}{m} \begin{bmatrix} 1 & -\frac{1}{N} & -\frac{1}{N} \\ -\frac{1}{N} & 1 & -\frac{1}{N} \\ -\frac{1}{N} & -\frac{1}{N} & 1 \end{bmatrix}$$



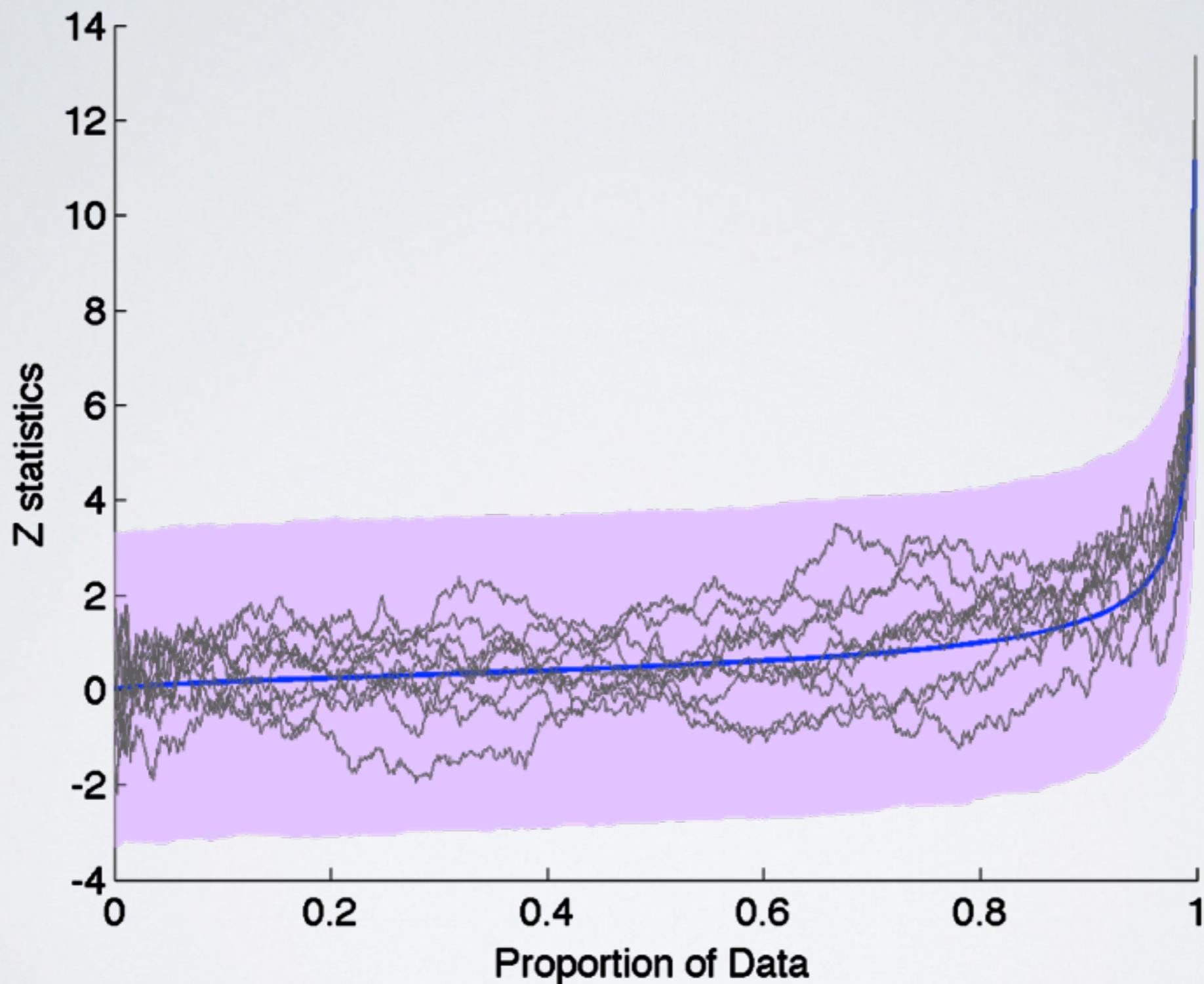
Random Walk of t_j

- Two assumptions about Y_j , when m is large $Y_1 \subseteq \mathcal{Y} \rightarrow Y_2 \subseteq \mathcal{Y} \rightarrow \dots$
 - $\{\bar{y}_j\}$ is jointly Normal
 - Accurate estimate of σ_y , $t_j = \frac{\overline{\{\bar{y}_j\}} - \mu_0}{\hat{\sigma}_j} \approx \frac{\overline{\{\bar{y}_j\}} - \mu_0}{\sigma_j} \sim \mathcal{N}(t_j)$

$$P(t_j | t_{j-1}) = \mathcal{N}(\mu_t(t_{j-1}), \Sigma_t)$$

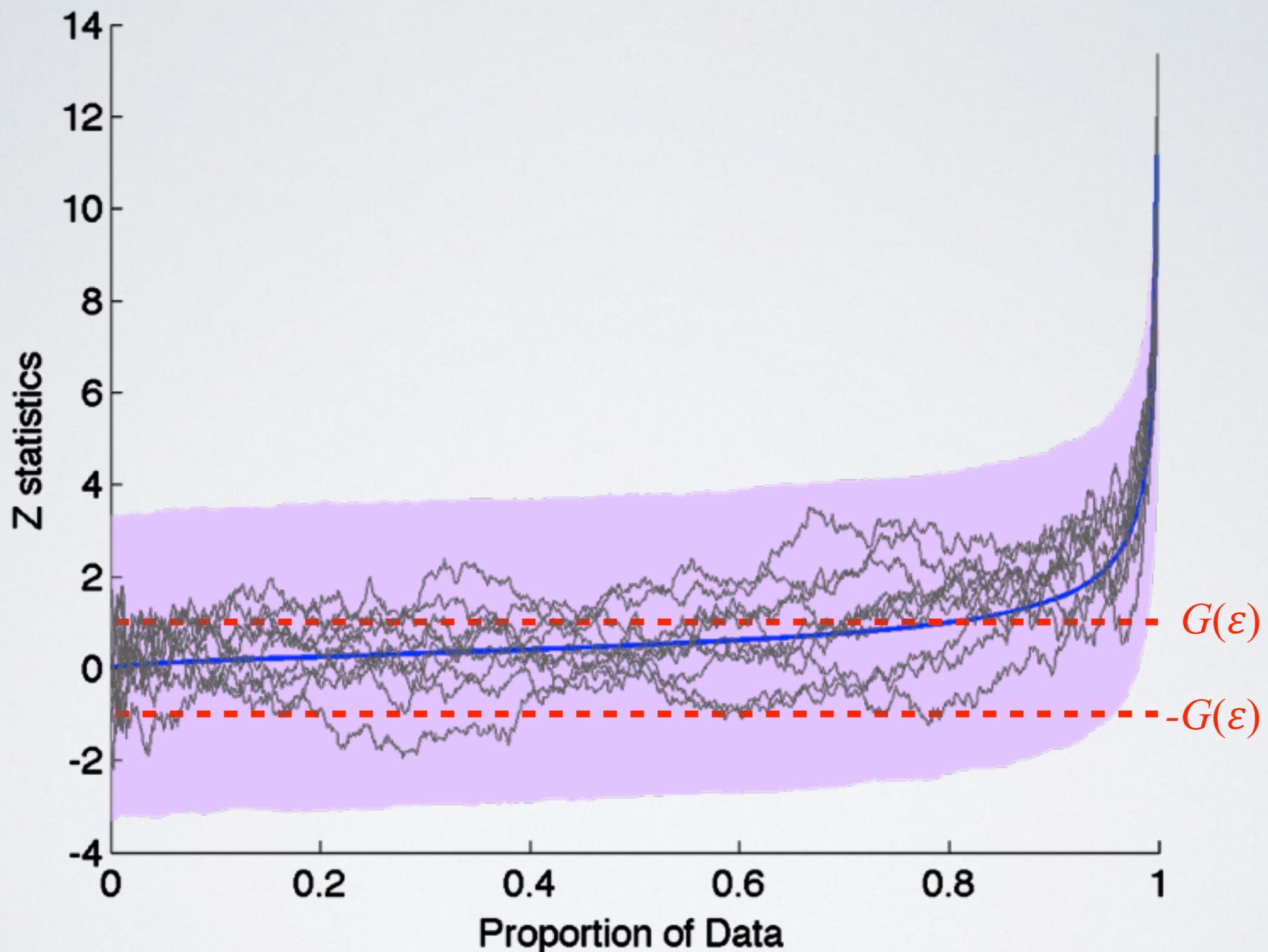


Example: $\mu > \mu_0$



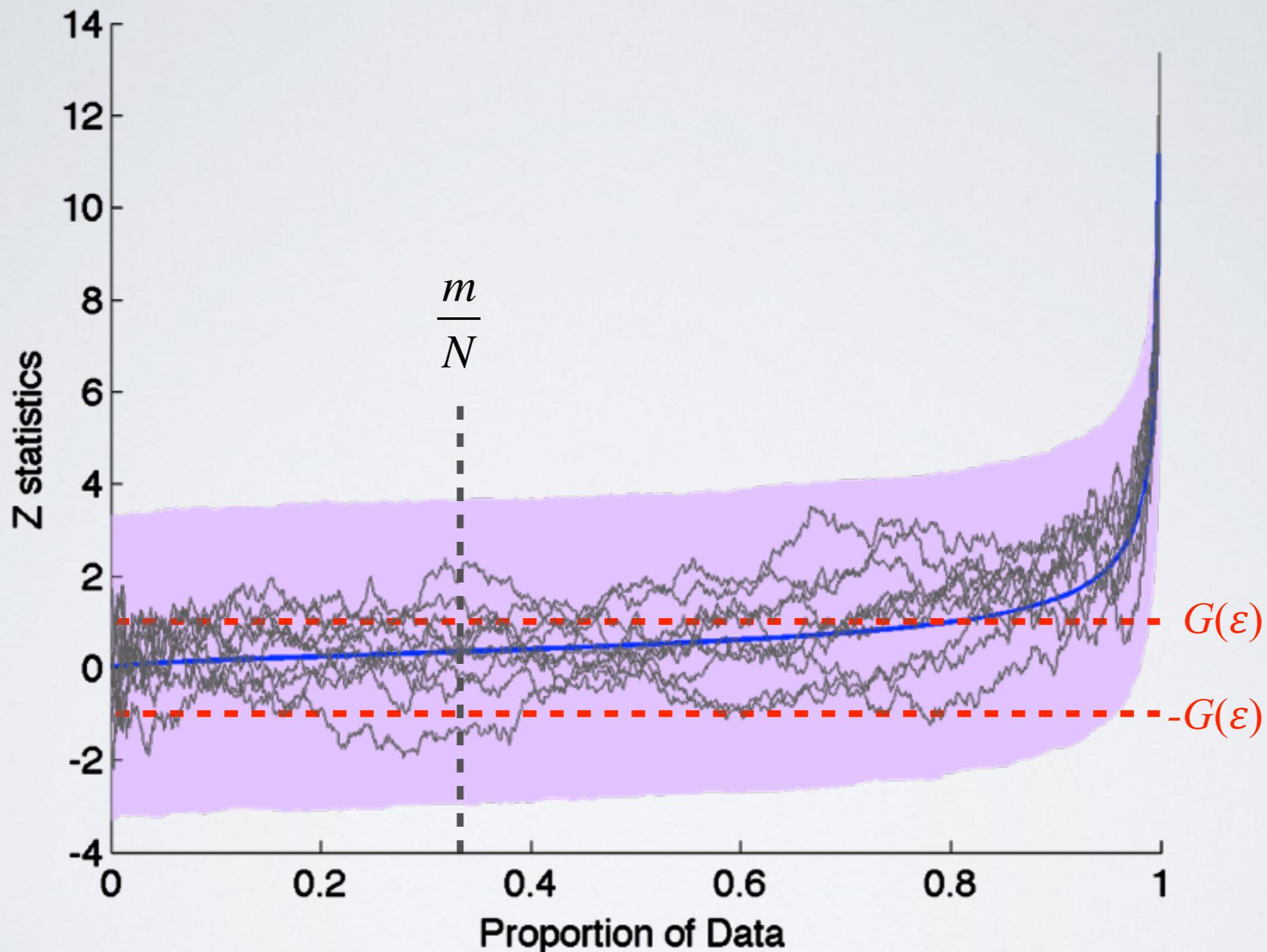


Example: $\mu > \mu_0$



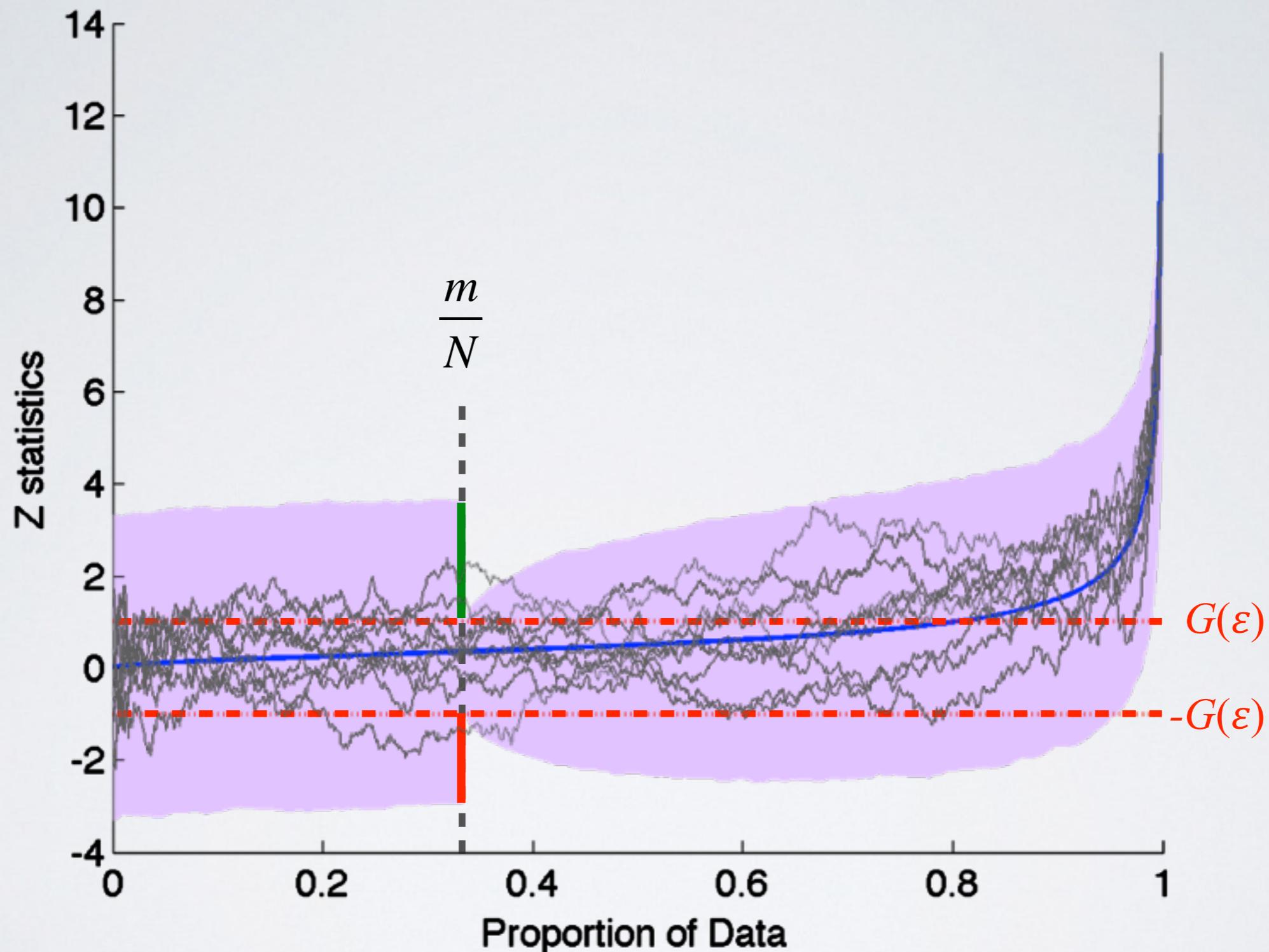


Example: $\mu > \mu_0$



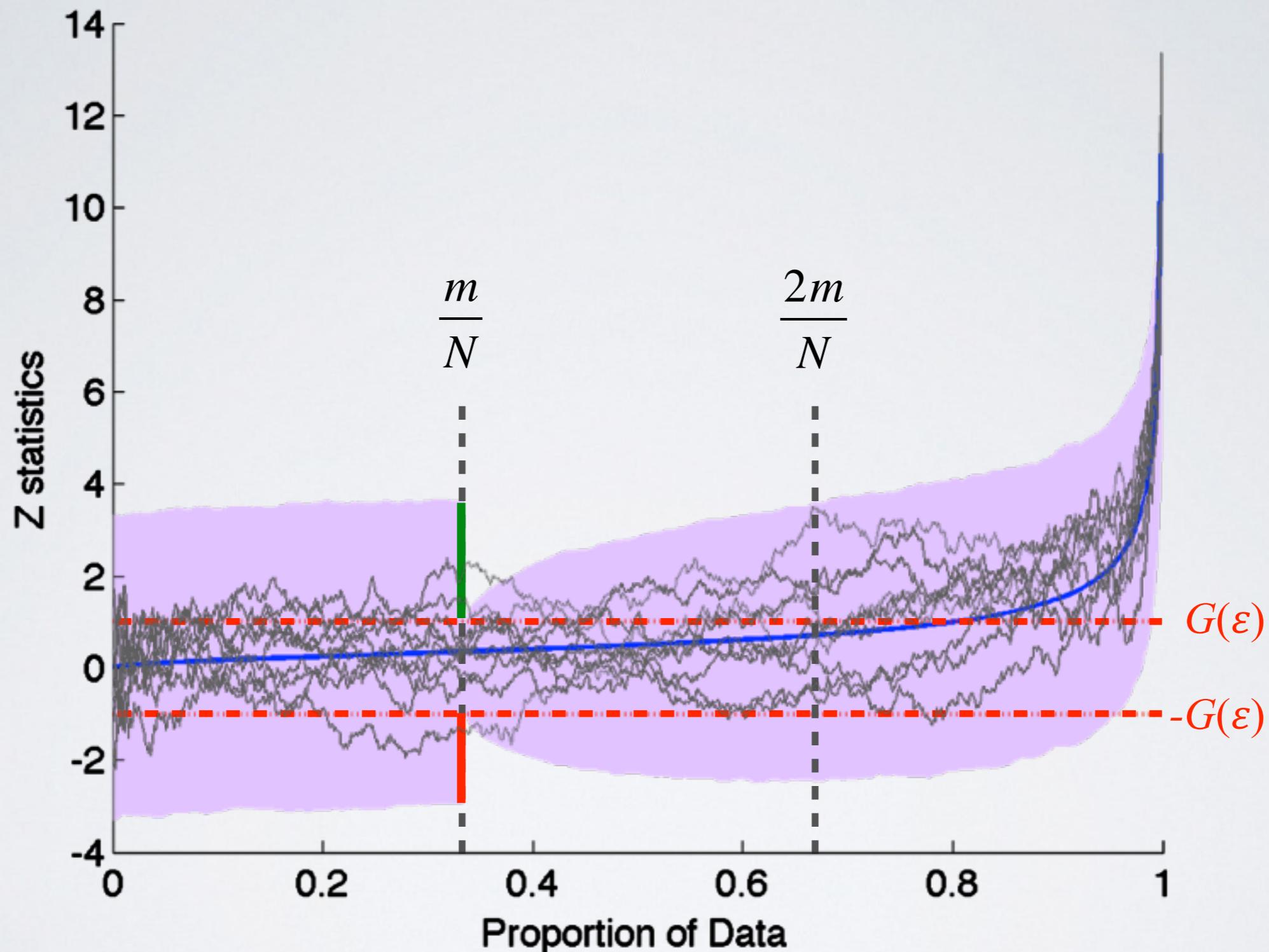


Example: $\mu > \mu_0$



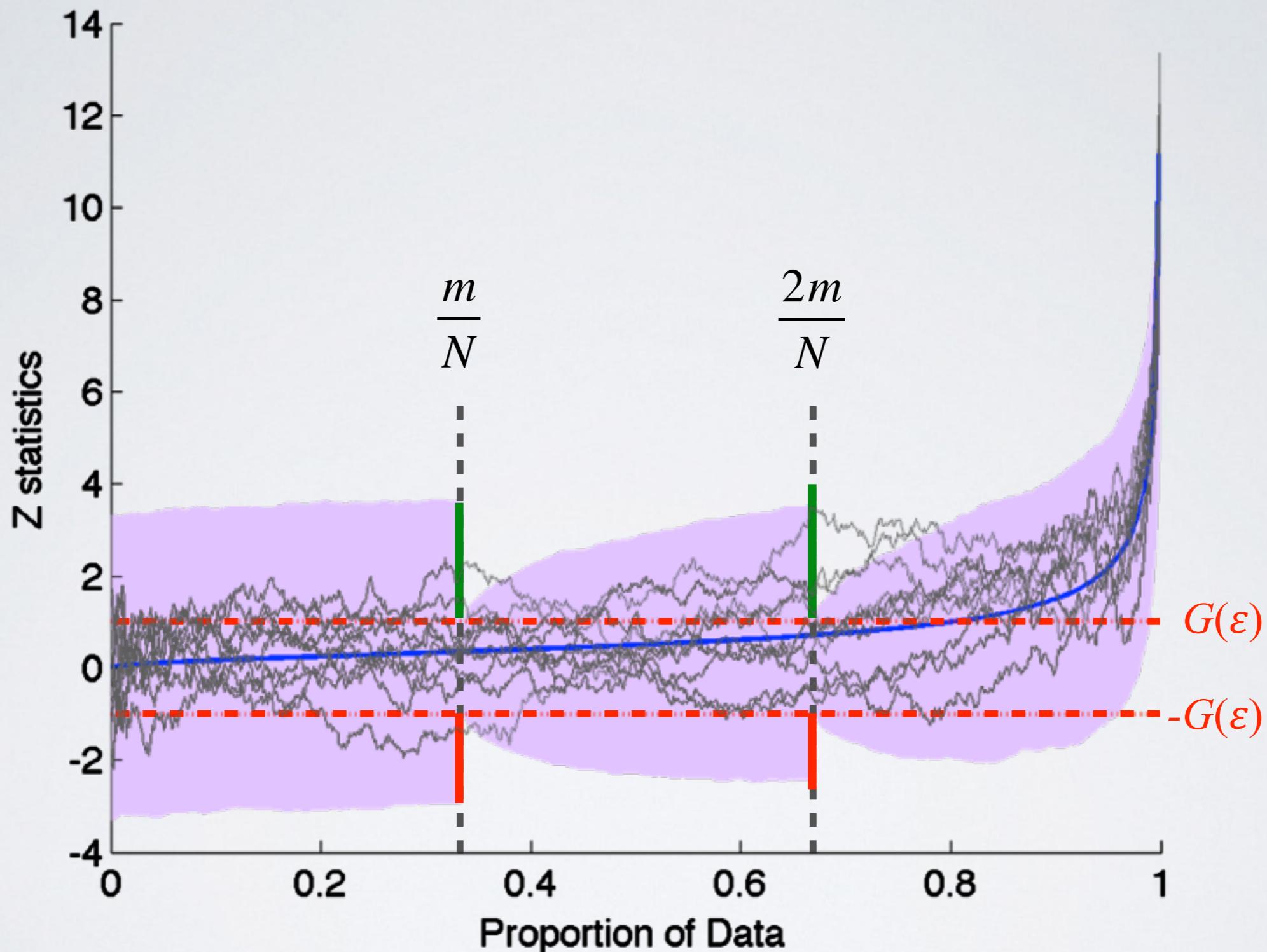


Example: $\mu > \mu_0$



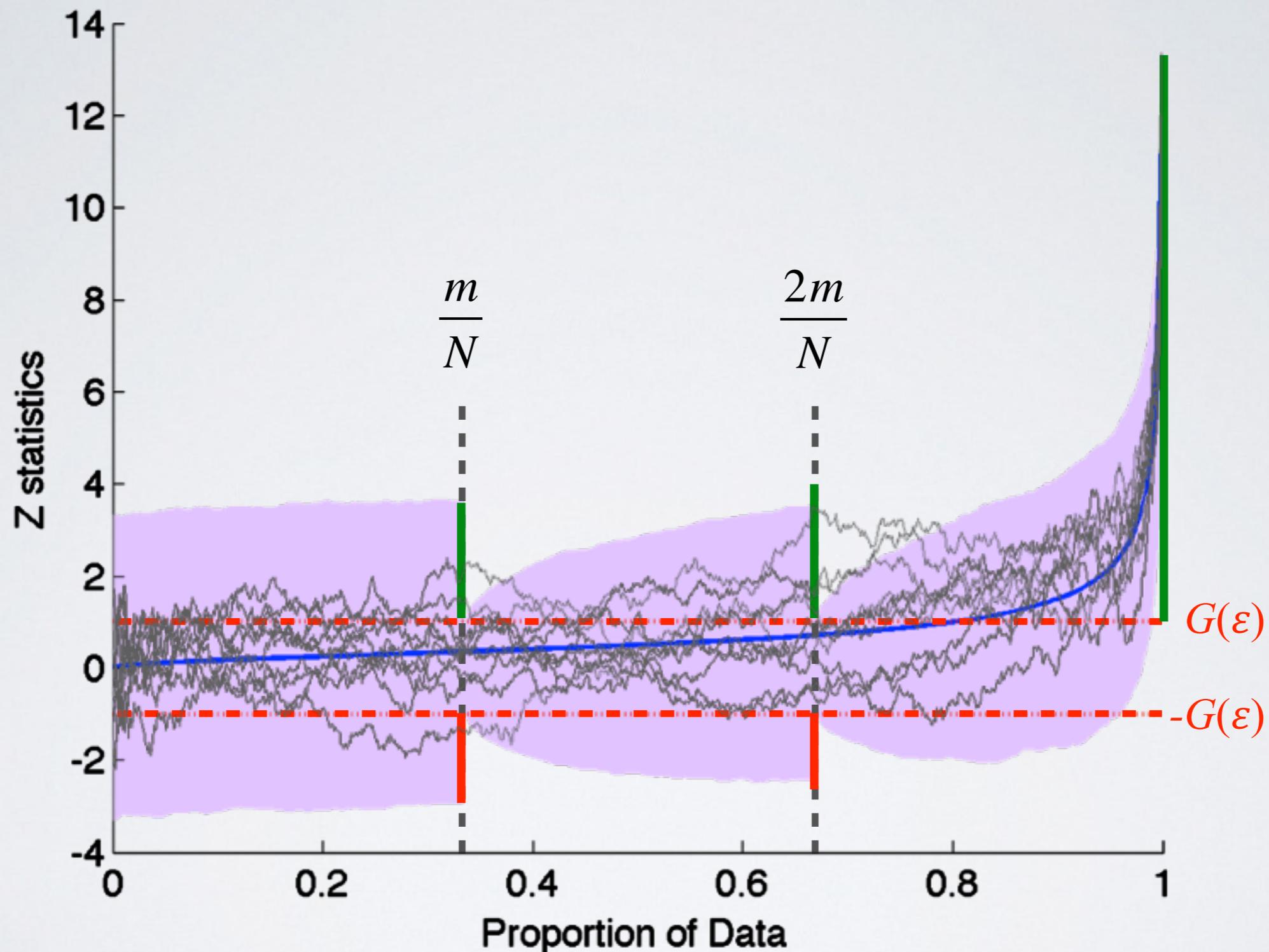


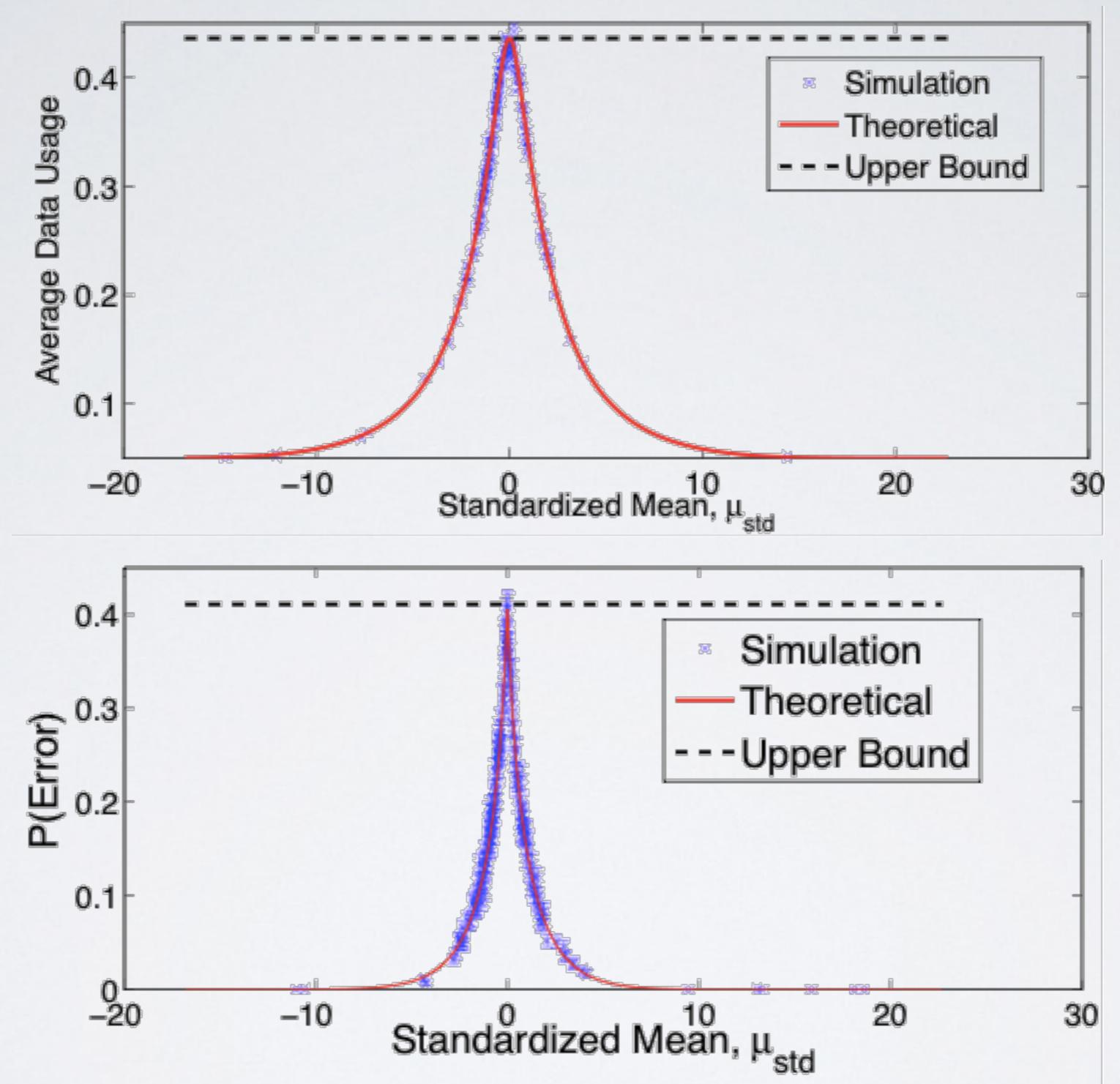
Example: $\mu > \mu_0$





Example: $\mu > \mu_0$

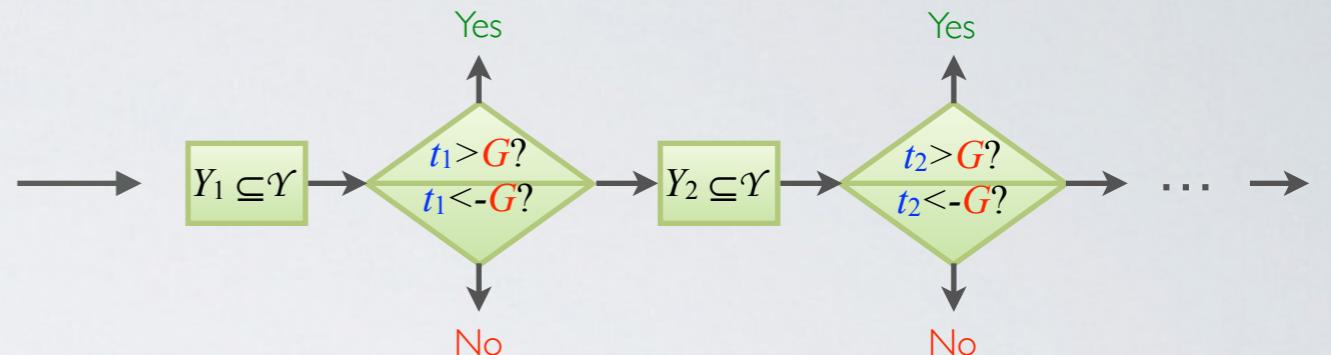




One Accept/Reject Step of M-H Iteration

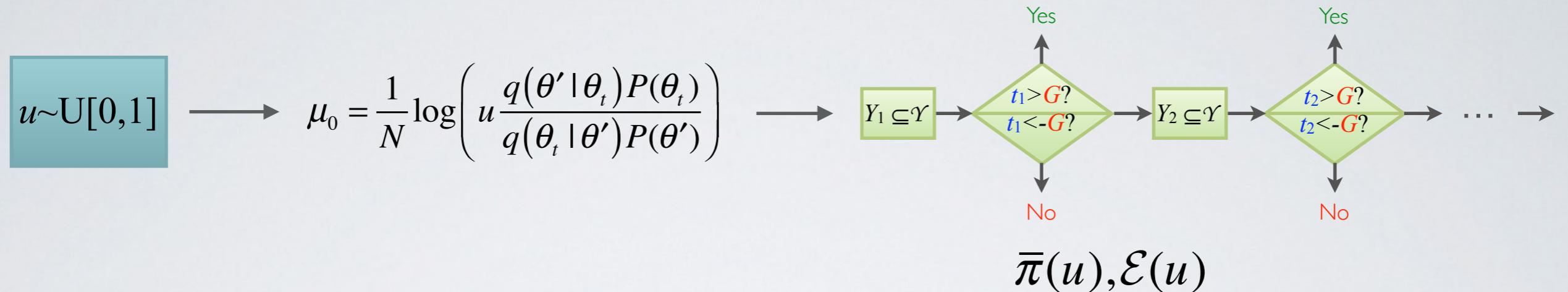
$$u \sim U[0,1]$$

$$\mu_0 = \frac{1}{N} \log \left(u \frac{q(\theta' | \theta_t) P(\theta_t)}{q(\theta_t | \theta') P(\theta')} \right)$$



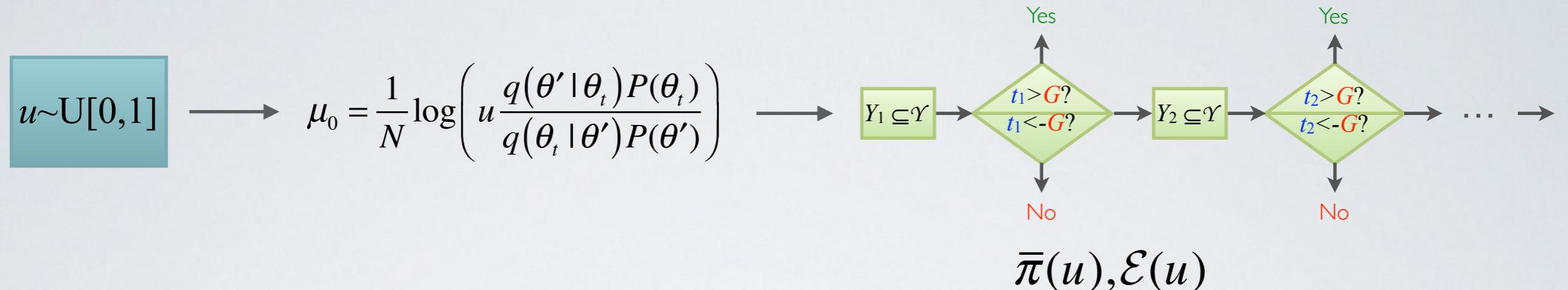
$$\bar{\pi}(u), \mathcal{E}(u)$$

One Accept/Reject Step of M-H Iteration



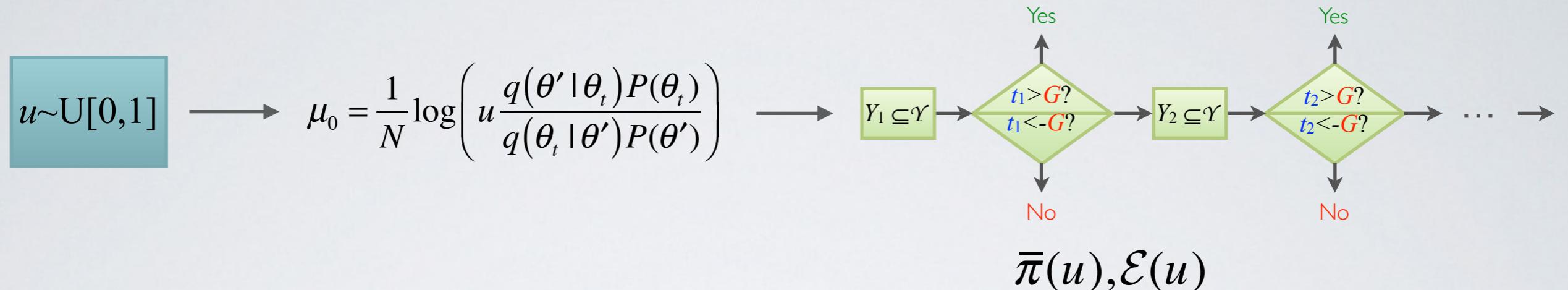
- Average data usage $\mathbb{E}[\bar{\pi}(u)] = \int_{u=0}^1 \bar{\pi}(u) du$

One Accept/Reject Step of M-H Iteration



- Average data usage $\mathbb{E}[\bar{\pi}(u)] = \int_{u=0}^1 \bar{\pi}(u) du$
- Error in the approximate acceptance probability

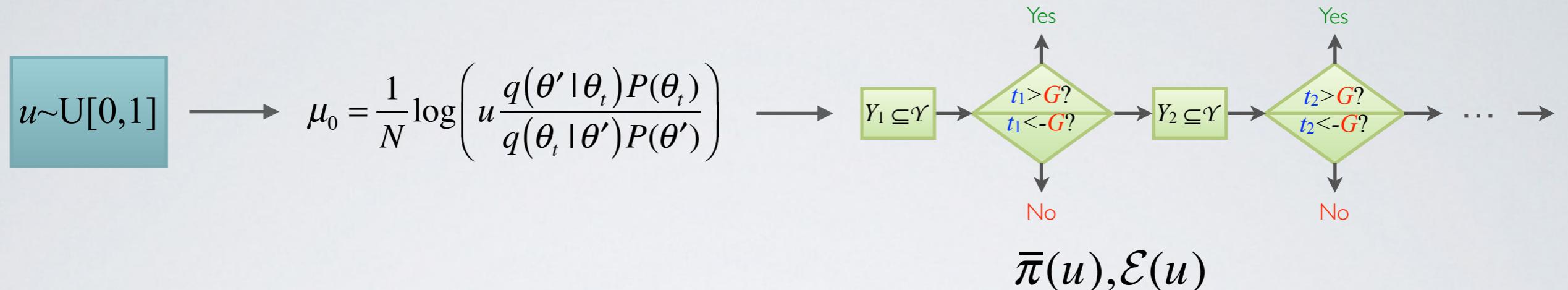
One Accept/Reject Step of M-H Iteration



- Average data usage $\mathbb{E}[\bar{\pi}(u)] = \int_{u=0}^1 \bar{\pi}(u) du$
- Error in the approximate acceptance probability

$$\Delta P_a = |P_{a,\epsilon} - P_a| = \left| \left[\int_0^{P_a} (1 - \mathcal{E}(u)) du + \int_{P_a}^1 \mathcal{E}(u) du \right] - \left[\int_0^{P_a} 1 du + \int_{P_a}^1 0 du \right] \right|$$

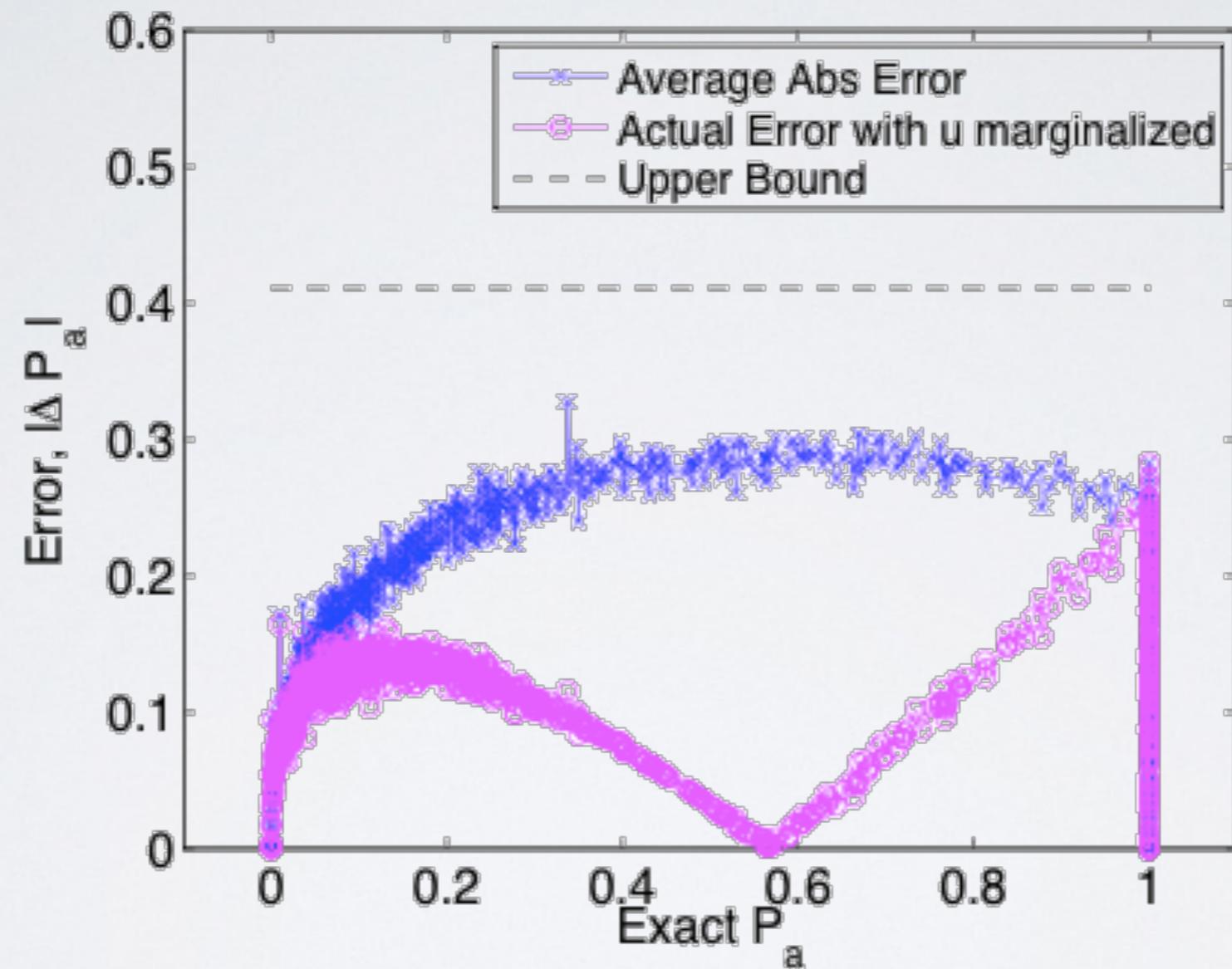
One Accept/Reject Step of M-H Iteration



- Average data usage $\mathbb{E}[\bar{\pi}(u)] = \int_{u=0}^1 \bar{\pi}(u) du$
- Error in the approximate acceptance probability

$$\begin{aligned}\Delta P_a &= |P_{a,\epsilon} - P_a| = \left| \left[\int_0^{P_a} (1 - \mathcal{E}(u)) du + \int_{P_a}^1 \mathcal{E}(u) du \right] - \left[\int_0^{P_a} 1 du + \int_{P_a}^1 0 du \right] \right| \\ &= \left| - \int_0^{P_a} \mathcal{E}(u) du + \int_{P_a}^1 \mathcal{E}(u) du \right|\end{aligned}$$

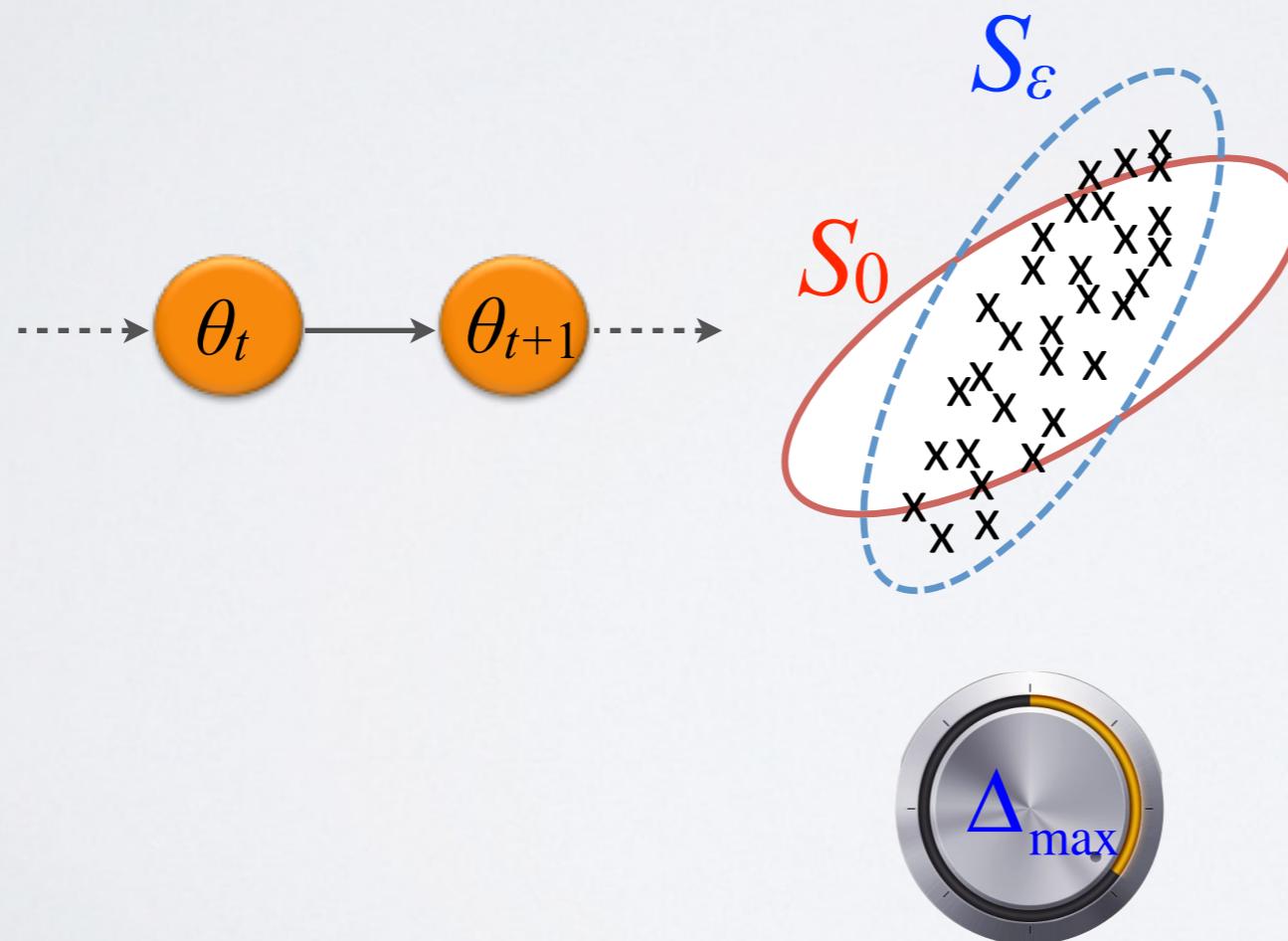
False negative False positive



Approximate stationary distribution, S_ε

Theorem 1. If the exact transition kernel \mathcal{T}_0 satisfies the contraction condition: $\exists \eta \in [0, 1)$ s.t. $d_v(P\mathcal{T}_0, S_0) \leq \eta d_v(P, S_0) \forall P$, the distance between the posterior distribution S_0 and the stationary distribution of our approximate Markov chain S_ε is upper bounded as:

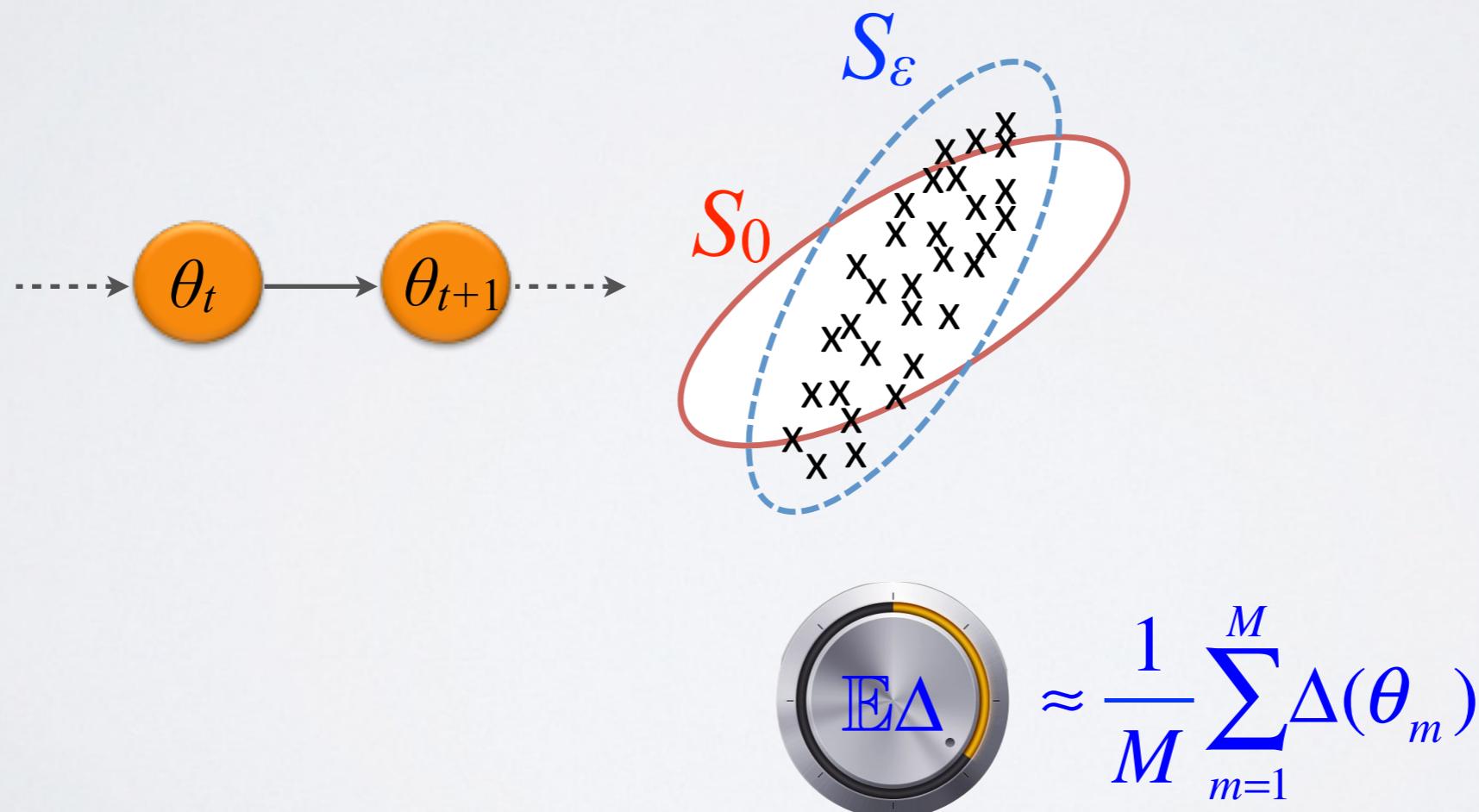
$$d_v(S_0, S_\varepsilon) \leq \frac{\Delta_{\max}}{1 - \eta}$$

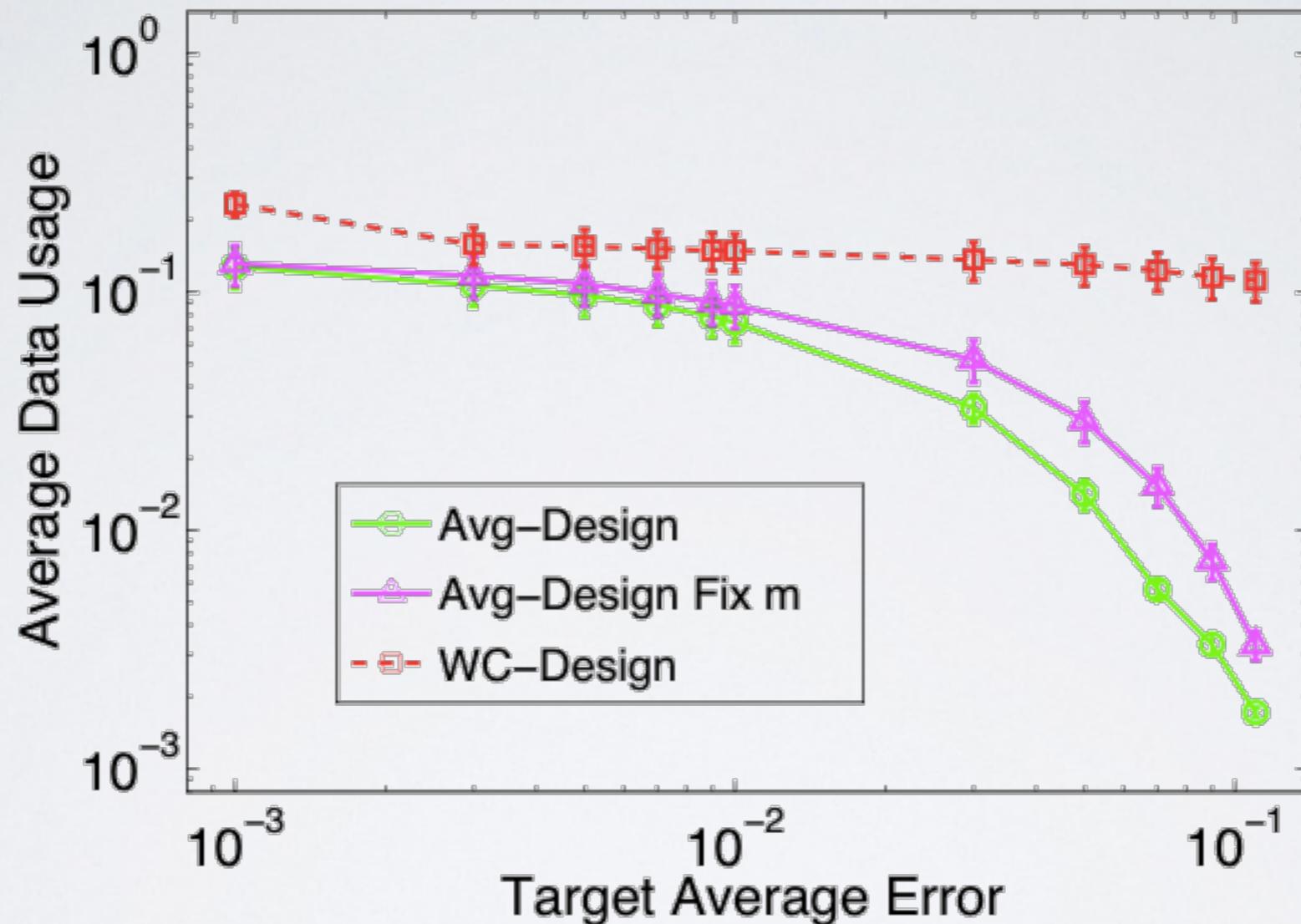


Approximate stationary distribution, S_ε

Theorem 1. If the exact transition kernel \mathcal{T}_0 satisfies the contraction condition: $\exists \eta \in [0, 1)$ s.t. $d_v(P\mathcal{T}_0, S_0) \leq \eta d_v(P, S_0) \forall P$, the distance between the posterior distribution S_0 and the stationary distribution of our approximate Markov chain S_ε is upper bounded as:

$$d_v(S_0, S_\varepsilon) \leq \frac{\Delta_{\max}}{1 - \eta}$$





Summary

- Proposed to design MCMC algorithms by balancing bias and variance under the big data setting.
- Proposed a scalable approximate M-H algorithm based on sequential hypothesis test on the accept/reject step.
- Analysis on the approximation error and speed gain.

Thanks!

Q&A

References

- Ahn, Sungjin, Anoop Korattikara, and Max Welling. "Bayesian posterior sampling via stochastic gradient Fisher scoring." In International Conference on Machine Learning, 2012.
- Bottou, Léon, and Olivier Bousquet. "The Tradeoffs of Large Scale Learning." NIPS. Vol. 4. 2007.
- Hoffman, Matthew D., et al. "Stochastic variational inference." The Journal of Machine Learning Research 14.1 (2013): 1303-1347.
- Singh, Sameer, Wick, Michael, and McCallum, Andrew. Monte Carlo MCMC: efficient inference by approximate sampling. In Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, pp. 1104–1113. Association for Computational Linguistics, 2012.
- Welling, Max, and Yee W. Teh. "Bayesian learning via stochastic gradient Langevin dynamics." Proceedings of the 28th International Conference on Machine Learning. 2011.