







Scaling Bayesian inference by constructing approximating exponential families

Jonathan H. Huggins

Postdoctoral Research Fellow Department of Biostatistics, Harvard

With: Ryan Adams, Tamara Broderick, James Zou

Microcredit



Microcredit



Cancer genomics





Microcredit



Cancer genomics







Fuel consumption

Microcredit



• Challenge: existing methods slow (and/or tedious, unreliable)

Microcredit



- Challenge: existing methods slow (and/or tedious, unreliable)
- Our proposal: use **efficient summaries** of data

Microcredit



- Challenge: existing methods slow (and/or tedious, unreliable)
- Our proposal: use **efficient summaries** of data
- Approximate sufficient statistics for simple, scalable
 Bayesian inference with error bounds for finite data

Roadmap

- Approximate Bayes review
- Likelihood approximation and dataset compression
- Approximate sufficient statistics
- Accuracy guarantees

Roadmap

- Approximate Bayes review
- Likelihood approximation and dataset compression
- Approximate sufficient statistics
- Accuracy guarantees

 $\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$



$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$$

• calculate...



$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$$

- calculate...
 - mean($\theta_{20} \mid Y$)



$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$$

- calculate...
 - mean($\theta_{20} \mid Y$)
 - var(θ₂₀ | Y)



$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$$

- calculate...
 - mean($\theta_{20} \mid Y$)
 - $var(\theta_{20} \mid Y)$
 - $\Pr[\theta_{20} < .1 \mid Y]$



$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$$

- calculate...
 - mean($\theta_{20} \mid Y$)
 - var(θ₂₀ | Y)
 - $\Pr[\theta_{20} < .1 \mid Y]$
- Goal: compute expectations wrt π



$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$$

- calculate...
 - mean($\theta_{20} \mid Y$)
 - var(θ₂₀ | Y)
 - $\Pr[\theta_{20} < .1 \mid Y]$
- Goal: compute expectations wrt π
- A hard problem!



$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$$

- calculate...
 - mean($\theta_{20} \mid Y$)
 - var(θ₂₀ | Y)
 - $\Pr[\theta_{20} < .1 \mid Y]$
- Goal: compute expectations wrt π
- A hard problem!



$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$$

- calculate...
 - mean($\theta_{20} \mid Y$)
 - var(θ₂₀ | Y)
 - $\Pr[\theta_{20} < .1 \mid Y]$
- Goal: compute expectations wrt π
- A hard problem!



$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$$

- calculate...
 - mean($\theta_{20} \mid Y$)
 - var(θ₂₀ | Y)
 - $\Pr[\theta_{20} < .1 \mid Y]$
- Goal: compute expectations wrt π
- A hard problem!



$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$$

- calculate...
 - mean($\theta_{20} \mid Y$)
 - var(θ₂₀ | Y)
 - $\Pr[\theta_{20} < .1 \mid Y]$
- Goal: compute expectations wrt π
- A hard problem!
- **Solution:** approximate π



1. Scalability

- 1. Scalability
 - large datasets

- 1. Scalability
 - large datasets
 - streaming and distributed data

- 1. Scalability
 - large datasets
 - streaming and distributed data
 - moderate-sized data with complex models

- 1. Scalability
 - large datasets
 - streaming and distributed data
 - moderate-sized data with complex models
- 2. Arbitrary accuracy: $d(\pi, \pi_{approx}) \rightarrow 0$

- 1. Scalability
 - large datasets
 - streaming and distributed data
 - moderate-sized data with complex models
- 2. Arbitrary accuracy: $d(\pi, \pi_{approx}) \rightarrow 0$
- 3. Validation of approximation quality: $d(\pi, \pi_{approx}) < c$

1. Scalability

- 2. Arbitrary accuracy
- 3. Validation of approximation quality

1. Scalability

Markov chain Monte Carlo

- 2. Arbitrary accuracy
- 3. Validation of approximation quality

1. Scalability

Markov chain Monte Carlo

- 2. Arbitrary accuracy
- 3. Validation of approximation quality

1. Scalability

Markov chain Monte Carlo

subsampling MCMC

- 2. Arbitrary accuracy
- 3. Validation of approximation quality



3. Validation of approximation quality

1. Scalability

- 2. Arbitrary accuracy
- 3. Validation of approximation quality

Markov chain Monte Carlo

subsampling MCMC

variational Bayes



3. Validation of approximation quality
Modern Bayesian inference

1. Scalability

- 2. Arbitrary accuracy
- 3. Validation of approximation quality

Markov chain Monte Carlo

subsampling MCMC

variational Bayes

consensus methods

Modern Bayesian inference

1. Scalability,

- 2. Arbitrary accuracy
- 3. Validation of approximation quality

Markov chain Monte Carlo

subsampling MCMC

variational Bayes

consensus methods

Modern Bayesian inference

1. Scalability

- 2. Arbitrary accuracy
- 3. Validation of approximation quality

Markov chain Monte Carlo

subsampling MCMC

variational Bayes

consensus methods

this talk: we get all three!

Roadmap

- Approximate Bayes review
- Likelihood approximation and dataset compression
- Approximate sufficient statistics
- Accuracy guarantees

Monte Carlo

 $\mathbb{E}[f(\theta) \mid Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t)$ Monte Carlo

Monte Carlo $\mathbb{E}[f(\theta) \mid Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \stackrel{\text{i.i.d.}}{\sim} \pi(\theta \mid Y)$

Monte Carlo $\mathbb{E}[f(\theta) \mid Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta \mid Y)$

Markov chain Monte Carlo (MCMC)

Monte Carlo $\mathbb{E}[f(\theta) | Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \stackrel{\text{i.i.d.}}{\sim} \pi(\theta | Y)$ Markov chain Monte Carlo (MCMC) $\theta_t \text{ not independent}$

Monte Carlo $\mathbb{E}[f(\theta) | Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \stackrel{\text{i.i.d.}}{\sim} \pi(\theta | Y)$ Markov chain Monte Carlo (MCMC) $\theta_t \text{ not independent}$



Monte Carlo $\mathbb{E}[f(\theta) | Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta | Y)$

Markov chain Monte Carlo (MCMC)

 θ_t not independent



Monte Carlo $\mathbb{E}[f(\theta) | Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta | Y)$

Markov chain Monte Carlo (MCMC)

 $\theta_t not$ independent

evaluate $\log p(Y \mid \theta) \pi_0(\theta)$ at θ_0 and θ'_1



Monte Carlo $\mathbb{E}[f(\theta) | Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta | Y)$

Markov chain Monte Carlo (MCMC)

 $\theta_t not$ independent

evaluate $\log p(Y \mid \theta) \pi_0(\theta)$ at θ_0 and θ'_1



 $\begin{array}{ll} \text{Monte Carlo} & \mathbb{E}[f(\theta) \mid Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) & \theta_t \stackrel{\text{i.i.d.}}{\sim} \pi(\theta \mid Y) \\ \\ \text{Markov chain Monte Carlo (MCMC)} & \theta_t & not \text{ independent} \end{array}$

 θ_0 θ_1 θ_0 θ_0'

Monte Carlo $\mathbb{E}[f(\theta) | Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \stackrel{\text{i.i.d.}}{\sim} \pi(\theta | Y)$ Markov chain Monte Carlo (MCMC) $\theta_t \text{ not independent}$



 $\mathbb{E}[f(\theta) \mid Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \stackrel{\text{i.i.d.}}{\sim} \pi(\theta \mid Y)$ Monte Carlo θ_t not independent

Markov chain Monte Carlo (MCMC)

 θ_1, θ_2

 θ_0

evaluate $\log p(Y \mid \theta) \pi_0(\theta)$ at θ_1 and θ'_2

Monte Carlo $\mathbb{E}[f(\theta) | Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta | Y)$

Markov chain Monte Carlo (MCMC)

 θ_t not independent



Monte Carlo
$$\mathbb{E}[f(\theta) | Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta | Y)$$

Markov chain Monte Carlo (MCMC)

 θ_t not independent



evaluate $\log p(Y \mid \theta) \pi_0(\theta)$ at θ_2 and θ'_3

Monte Carlo $\mathbb{E}[f(\theta) | Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta | Y)$

Markov chain Monte Carlo (MCMC)

 θ_t not independent



evaluate $\log p(Y \mid \theta) \pi_0(\theta)$ at θ_2 and θ'_3

Monte Carlo $\mathbb{E}[f(\theta) | Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta | Y)$

Markov chain Monte Carlo (MCMC)

 θ_t not independent



 $\begin{array}{ll} \text{Monte Carlo} & \mathbb{E}[f(\theta) \mid Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) & \theta_t \stackrel{\text{i.i.d.}}{\sim} \pi(\theta \mid Y) \\ \text{Markov chain Monte Carlo (MCMC)} & \theta_t \text{ not independent} \\ & Y = \{y_1, y_2, \dots, y_N\} \\ & \bullet_{\theta_0} \end{array}$



 θ_T

 $\mathbb{E}[f(\theta) \mid Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \stackrel{\text{i.i.d.}}{\sim} \pi(\theta \mid Y)$ **Monte Carlo** θ_t not independent Markov chain Monte Carlo (MCMC) $Y = \{y_1, y_2, \dots, y_N\}$ θ_1, θ_2 **Problem:** $\log p(Y|\theta)$ is expensive to evaluate if N is large θ_0

 $\mathbb{E}[f(\theta) \mid Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \stackrel{\text{i.i.d.}}{\sim} \pi(\theta \mid Y)$ **Monte Carlo** θ_t not independent Markov chain Monte Carlo (MCMC) $Y = \{y_1, y_2, \dots, y_N\}$ θ_1, θ_2 **Problem:** $\log p(Y|\theta)$ is expensive $\Omega(N \times T)$ to evaluate if N is large time θ_0 $heta_T$ θ_3

Monte Carlo
$$\mathbb{E}[f(\theta) | Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta | Y)$$

Markov chain Monte Carlo (MCMC)

 θ_t not independent

$$Y = \{y_1, y_2, \dots, y_N\}$$



Problem: $\log p(Y|\theta)$ is expensive to evaluate if *N* is large

 $\Omega(N \times T)$ time

Solution:

Monte Carlo
$$\mathbb{E}[f(\theta) \mid Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta \mid Y)$$

Markov chain Monte Carlo (MCMC)

 θ_t not independent

$$Y = \{y_1, y_2, \dots, y_N\}$$



Problem: $\log p(Y|\theta)$ is expensive to evaluate if *N* is large

 $\Omega(N imes T)$ time

Solution:

1. replace $\log p(Y|\theta)$ with a fast-tocompute proxy $\ell(\theta, g(Y))$

Monte Carlo
$$\mathbb{E}[f(\theta) \mid Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta \mid Y)$$

Markov chain Monte Carlo (MCMC)

 θ_t not independent

$$Y = \{y_1, y_2, \dots, y_N\}$$



Problem: $\log p(Y|\theta)$ is expensive to evaluate if *N* is large

 $\Omega(N imes T)$ time

Solution:

1. replace $\log p(Y|\theta)$ with a fast-tocompute proxy $\ell(\theta, g(Y))$ short summary of Y

Monte Carlo
$$\mathbb{E}[f(\theta) \mid Y] \approx T^{-1} \sum_{t=1}^{T} f(\theta_t) \quad \theta_t \overset{\text{i.i.d.}}{\sim} \pi(\theta \mid Y)$$

Markov chain Monte Carlo (MCMC)

 θ_t not independent

$$Y = \{y_1, y_2, \dots, y_N\}$$



Problem: $\log p(Y|\theta)$ is expensive to evaluate if *N* is large

 $\Omega(N imes T)$ time

Solution:

- 1. replace $\log p(Y|\theta)$ with a fast-tocompute proxy $\ell(\theta, g(Y))$
- 2. choose $\ell(\theta, g(Y))$ so approximate posterior is accurate

 $Y = \{y_1, y_2, \dots, y_N\}$

 $Y = \{y_1, y_2, \dots, y_N\}$

 $\log p(y_n|\theta) = \eta(\theta) \cdot \tau(y_n)$

 $Y = \{y_1, y_2, \dots, y_N\}$ $\log p(y_n | \theta) = \eta(\theta) \cdot \tau(y_n)$ \checkmark reparameterization

 $Y = \{y_1, y_2, \dots, y_N\}$ $\log p(y_n | \theta) = \eta(\theta) \cdot \tau(y_n)$ \checkmark reparameterization sufficient statistics

$$Y = \{y_1, y_2, \dots, y_N\}$$
$$\sum_{n=1}^N \log p(y_n | \theta) = \eta(\theta) \cdot \sum_{n=1}^N \tau(y_n)$$

 $Y = \{y_1, y_2, \dots, y_N\}$

 $\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \sum_{n=1}^{N} \tau(y_n)$

$$Y = \{y_1, y_2, \dots, y_N\}$$

$$\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \underbrace{\sum_{n=1}^{N} \tau(y_n)}_{g(Y)}$$

$$Y = \{y_1, y_2, \dots, y_N\}$$
$$\log p(Y|\theta) = \sum_{n=1}^N \log p(y_n|\theta) = \eta(\theta) \cdot \underbrace{\sum_{n=1}^N \tau(y_n)}_{g(Y)}$$

Run MCMC for *T* iterations:
$$Y = \{y_1, y_2, \dots, y_N\}$$

$$\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \underbrace{\sum_{n=1}^{N} \tau(y_n)}_{g(Y)}$$

Run MCMC for *T* iterations:

• Naively: $\Omega(N \times T)$ time

$$Y = \{y_1, y_2, \dots, y_N\}$$

$$\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \underbrace{\sum_{n=1}^{N} \tau(y_n)}_{g(Y)}$$

Run MCMC for T iterations: $\log p(Y|\theta)$ takes $\Omega(N)$ time to evaluateNaively: $\Omega(N \times T)$ time

$$Y = \{y_1, y_2, \dots, y_N\}$$

$$\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \underbrace{\sum_{n=1}^{N} \tau(y_n)}_{g(Y)}$$

Run MCMC for T iterations: $\log p(Y|\theta)$ takes $\Omega(N)$ • Naively: $\Omega(N \times T)$ time

$$Y = \{y_1, y_2, \dots, y_N\}$$

$$\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \underbrace{\sum_{n=1}^{N} \tau(y_n)}_{g(Y)}$$

Run MCMC for T iterations:

- Naively: $\Omega(N \times T)$ time
- Using EF structure: O(N + T) time

 $\log p(Y|\theta)$ takes $\Omega(N)$

time to evaluate

$$Y = \{y_1, y_2, \dots, y_N\}$$

$$\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \underbrace{\sum_{n=1}^{N} \tau(y_n)}_{g(Y)}$$

Run MCMC for T iterations:

- Naively: $\Omega(N)$
- IS: $\log p(Y|\theta)$ takes $\Omega(N)$ time to evaluate $\Omega(N \times T)$ time
 - Using EF structure:

compute g(Y)

O(N+T) time

$$Y = \{y_1, y_2, \dots, y_N\}$$

$$\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \underbrace{\sum_{n=1}^{N} \tau(y_n)}_{g(Y)}$$

10

Run MCMC for *T* iterations:

- Naively: $\Omega(N \times T)$
- Using EF structure:

ations: $\log p(Y|\theta)$ takes $\Omega(N)$ time to evaluate $\Omega(N \times T)$ time O(N + T) time $\log p(Y|\theta)$ takes O(1)time to evaluate

$$Y = \{y_1, y_2, \dots, y_N\}$$

$$\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \underbrace{\sum_{n=1}^{N} \tau(y_n)}_{g(Y)}$$



STREAMING

STREAMING

 $\tau = 0$

STREAMING

 $\tau = 0$

for *n* = 1, ..., *N*

STREAMING

 $\tau = 0$ for n = 1, ..., N $\tau = \tau + \tau(y_n)$

STREAMING

τ = 0for n = 1, ..., N $τ = τ + τ(y_n)$ construct log p(Y|θ) using τ

STREAMING

 $\tau = 0$ for n = 1, ..., N $\tau = \tau + \tau(y_n)$ construct log $p(Y|\theta)$ using τ run MCMC

$$\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \sum_{n=1}^{N} \tau(y_n)$$

STREAMING

DISTRIBUTED

 $\tau = 0$

for *n* = 1, ..., *N*

 $\tau = \tau + \tau(y_n)$

construct $\log p(Y|\theta)$ using τ run MCMC

$$\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \sum_{n=1}^{N} \tau(y_n)$$

STREAMING

 $\tau = 0$

for n = 1, ..., N $\tau = \tau + \tau(y_n)$ construct $\log p(Y|\theta)$ using τ run MCMC **DISTRIBUTED** for b = 1, ..., B in parallel

 $\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \sum_{n=1}^{N} \tau(y_n)$

STREAMING

 $\tau = 0$

for *n* = 1, ..., *N*

 $\tau = \tau + \tau(y_n)$

construct $\log p(Y|\theta)$ using τ run MCMC

DISTRIBUTED

for b = 1, ..., B in parallel

1. worker *b* reads data subset $Y_b = \{ y_{B(b-1)/n}, \dots, y_{Bb/n-1} \}$

 $\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \sum_{n=1}^{N} \tau(y_n)$

STREAMING

 $\tau = 0$

for *n* = 1, ..., *N*

$$\tau = \tau + \tau(y_n)$$

construct $\log p(Y|\theta)$ using τ run MCMC

DISTRIBUTED

for b = 1, ..., B in parallel

- 1. worker *b* reads data subset $Y_b = \{ y_{B(b-1)/n}, \dots, y_{Bb/n-1} \}$
- 2. worker *b* computes $\tau_b = \tau(Y_b)$

 $\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \sum_{n=1}^{N} \tau(y_n)$

STREAMING

 $\tau = 0$

for *n* = 1, ..., *N*

 $\tau = \tau + \tau(y_n)$

construct $\log p(Y|\theta)$ using τ run MCMC

DISTRIBUTED

for b = 1, ..., B in parallel

- 1. worker *b* reads data subset $Y_b = \{ y_{B(b-1)/n}, \dots, y_{Bb/n-1} \}$
- 2. worker *b* computes $\tau_b = \tau(Y_b)$

$$\tau = \tau_1 + \cdots + \tau_b$$

 $\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \sum_{n=1}^{N} \tau(y_n)$

STREAMING

 $\tau = 0$

for *n* = 1, ..., *N*

 $\tau = \tau + \tau(y_n)$

construct $\log p(Y|\theta)$ using τ run MCMC

DISTRIBUTED

for b = 1, ..., B in parallel

- 1. worker *b* reads data subset $Y_b = \{ y_{B(b-1)/n}, \dots, y_{Bb/n-1} \}$
- 2. worker *b* computes $\tau_b = \tau(Y_b)$

$$\tau = \tau_1 + \cdots + \tau_b$$

construct $\log p(Y|\theta)$ using τ

 $\log p(Y|\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) = \eta(\theta) \cdot \sum_{n=1}^{N} \tau(y_n)$

STREAMING

 $\tau = 0$

for *n* = 1, ..., *N*

 $\tau = \tau + \tau(y_n)$

construct $\log p(Y|\theta)$ using τ run MCMC

DISTRIBUTED

for b = 1, ..., B in parallel

- 1. worker *b* reads data subset $Y_b = \{ y_{B(b-1)/n}, \dots, y_{Bb/n-1} \}$
- 2. worker *b* computes $\tau_b = \tau(Y_b)$

 $\tau = \tau_1 + \cdots + \tau_b$

construct $\log p(Y|\theta)$ using τ run MCMC

Streaming and distributed exponential family inference



Roadmap

- Approximate Bayes review
- Likelihood approximation and dataset compression
- Approximate sufficient statistics
- Accuracy guarantees

Polynomials are good for approximation

 $-5.8x^6 + 14.9x^4 - 10.1x^2 + 1.1$

1. Computationally convenient

- 1. Computationally convenient
- 2. Can approximate any smooth function

- 1. Computationally convenient
- 2. Can approximate any smooth function
- 3. Approximation properties are well-understood

- 1. Computationally convenient
- 2. Can approximate any smooth function
- 3. Approximation properties are well-understood

Polynomials are good for approximation

 $-5.8x^6 + 14.9x^4 - 10.1x^2 + 1.1$

- 1. Computationally convenient
- 2. Can approximate any smooth function
- 3. Approximation properties are well-understood



Polynomial approximate sufficient statistics (PASS) $Y = \{y_1, y_2, \dots, y_N\}$ Ĥ

 $Y = \{y_1, y_2, \dots, y_N\}$ $X = \{x_1, x_2, \dots, x_N\}$



$$Y = \{y_1, y_2, \dots, y_N\}$$
$$X = \{x_1, x_2, \dots, x_N\}$$
$$x_n = (x_{n1}, x_{n2}, \dots, x_{nd})$$



 $Y = \{y_1, y_2, \dots, y_N\}$ $X = \{x_1, x_2, \dots, x_N\}$ $x_n = (x_{n1}, x_{n2}, \dots, x_{nd})$ $\log p(y_n | x_n, \theta) \approx \eta(\theta) \cdot \tau(y_n, x_n)$



$$Y = \{y_1, y_2, \dots, y_N\}$$
$$X = \{x_1, x_2, \dots, x_N\}$$
$$x_n = (x_{n1}, x_{n2}, \dots, x_{nd})$$
$$\log p(y_n | x_n, \theta) \approx \eta(\theta) \cdot \tau(y_n, x_n)$$
$$\tau(y_n, x_n) = (y_n, x_{n1}, x_{n2}, \dots, x_{nd}, x_n)$$



$$Y = \{y_1, y_2, \dots, y_N\}$$
$$X = \{x_1, x_2, \dots, x_N\}$$
$$x_n = (x_{n1}, x_{n2}, \dots, x_{nd})$$
$$\log p(y_n | x_n, \theta) \approx \eta(\theta) \cdot \tau(y_n, x_n)$$
$$\tau(y_n, x_n) = (y_n, x_{n1}, x_{n2}, \dots, x_{nd},$$
$$y_n^2, x_{n1}^2, x_{n2}^2, \dots, x_{nd}^2,$$


$Y = \{y_1, y_2, \dots, y_N\}$ $X = \{x_1, x_2, \dots, x_N\}$ $x_n = (x_{n1}, x_{n2}, \dots, x_{nd})$ $\log p(y_n | x_n, \theta) \approx \eta(\theta) \cdot \tau(y_n, x_n)$ $\tau(y_n, x_n) = (y_n, x_{n1}, x_{n2}, \dots, x_{nd},$ $y_n^2, x_{n1}^2, x_{n2}^2, \ldots, x_{nd}^2,$ $y_n x_{n1}, \ldots, y_n x_{nd},$



 $Y = \{y_1, y_2, \dots, y_N\}$ $X = \{x_1, x_2, \dots, x_N\}$ $x_n = (x_{n1}, x_{n2}, \dots, x_{nd})$ $\log p(y_n | x_n, \theta) \approx \eta(\theta) \cdot \tau(y_n, x_n)$ $\tau(y_n, x_n) = (y_n, x_{n1}, x_{n2}, \dots, x_{nd},$ $y_n^2, x_{n1}^2, x_{n2}^2, \ldots, x_{nd}^2,$ $y_n x_{n1}, \ldots, y_n x_{nd},$ $x_{n1}x_{n2}, x_{n1}x_{n3}, \ldots,$



 $Y = \{y_1, y_2, \dots, y_N\}$ $X = \{x_1, x_2, \dots, x_N\}$ $x_n = (x_{n1}, x_{n2}, \dots, x_{nd})$ $\log p(y_n | x_n, \theta) \approx \eta(\theta) \cdot \tau(y_n, x_n)$ $\tau(y_n, x_n) = (y_n, x_{n1}, x_{n2}, \dots, x_{nd},$ $y_n^2, x_{n1}^2, x_{n2}^2, \ldots, x_{nd}^2,$ $y_n x_{n1}, \ldots, y_n x_{nd},$ $x_{n1}x_{n2}, x_{n1}x_{n3}, \ldots,$



 $Y = \{y_1, y_2, \dots, y_N\}$ $X = \{x_1, x_2, \dots, x_N\}$ $x_n = (x_{n1}, x_{n2}, \dots, x_{nd})$ $\log p(y_n | x_n, \theta) \approx \eta(\theta) \cdot \tau(y_n, x_n)$ $\tau(y_n, x_n) = (y_n, x_{n1}, x_{n2}, \dots, x_{nd},$ $y_n^2, x_{n1}^2, x_{n2}^2, \ldots, x_{nd}^2,$ $y_n x_{n1}, \ldots, y_n x_{nd},$ $x_{n1}x_{n2}, x_{n1}x_{n3}, \ldots,$

 $y_n x_{n1}, \dots, y_n x_{nd},$ $x_{n1} x_{n2}, x_{n1} x_{n3}, \dots,$ $\dots,$ $y_n^M, x_{n1}^M, x_{n2}^M, \dots, x_{nd}^M)$





generalized linear models: $\log p(y_n | x_n, \theta) = \phi(y_n, \theta \cdot x_n)$

• Poisson regression (count data):

- Poisson regression (count data):
 - ➡ # of foreclosures by region



- Poisson regression (count data):
 - ➡ # of foreclosures by region
- Logistic regression (binary data)



- Poisson regression (count data):
 - ➡ # of foreclosures by region
- Logistic regression (binary data)
 - ➡ patient has cancer?





- Poisson regression (count data):
 - ➡ # of foreclosures by region
- Logistic regression (binary data)
 - ➡ patient has cancer?
- Robust regression (continuous data)





- Poisson regression (count data):
 - ➡ # of foreclosures by region
- Logistic regression (binary data)
 - ➡ patient has cancer?
- Robust regression (continuous data)
 - ➡ birth rate by region



- Poisson regression (count data):
 - ➡ # of foreclosures by region
- Logistic regression (binary data)
 - ➡ patient has cancer?
- Robust regression (continuous data)
 - ➡ birth rate by region



$$\tau(y_n, x_n) = \left(a(k, M)y_n^{k_0} \prod_{i=1}^d x_{ni}^{k_i}\right)_{\substack{k \in \mathbb{N}^{d+1} \\ \sum_i k \leq M}} \quad \eta(\theta) = \left(\prod_{i=1}^d \theta_i^{k_i}\right)_{\substack{k \in \mathbb{N}^{d+1} \\ \sum_i k \leq M}}$$





data $(x_n, y_n) \in \mathbb{R}^d \times \{\pm 1\}$

data $(x_n, y_n) \in \mathbb{R}^d \times \{\pm 1\}$ parameter $\theta \in \mathbb{R}^d$

data $(x_n, y_n) \in \mathbb{R}^d \times \{\pm 1\}$ parameter $\theta \in \mathbb{R}^d$ log-likelihood $\log p(y_n | x_n, \theta) = -\log(1 + e^{-y_n x_n \cdot \theta})$

data $(x_n, y_n) \in \mathbb{R}^d \times \{\pm 1\}$ parameter $\theta \in \mathbb{R}^d$ log-likelihood $\log p(y_n | x_n, \theta) = -\log(1 + e^{-y_n x_n \cdot \theta})$ $= \phi(y_n x_n \cdot \theta)$

data $(x_n, y_n) \in \mathbb{R}^d \times \{\pm 1\}$ parameter $\theta \in \mathbb{R}^d$ log-likelihood $\log p(y_n | x_n, \theta) = -\log(1 + e^{-y_n x_n \cdot \theta})$ $= \phi(y_n x_n \cdot \theta) \approx \phi_2(y_n x_n \cdot \theta)$



PASS-GLM effective in distributed and streaming settings

PASS-GLM effective in distributed and streaming settings

Criteo advertising dataset

PASS-GLM effective in distributed and streaming settings

Criteo advertising dataset



• Distributed

PASS-GLM effective in distributed and streaming settings

Criteo advertising dataset



- Distributed
 - 6M observations with 1K covariates

PASS-GLM effective in distributed and streaming settings

Criteo advertising dataset



- Distributed
 - 6M observations with 1K covariates
 - 16 seconds using 22 cores

PASS-GLM effective in distributed and streaming settings



- Distributed
 - 6M observations with 1K covariates
 - 16 seconds using 22 cores

• Streaming

PASS-GLM effective in distributed and streaming settings



- Distributed
 - 6M observations with 1K covariates
 - 16 seconds using 22 cores

- Streaming
 - 40M observations with 20K covariates

PASS-GLM effective in distributed and streaming settings



- Distributed
 - 6M observations with 1K covariates
 - 16 seconds using 22 cores

- Streaming
 - 40M observations with 20K covariates
 - Competitive with SGD

- Webspam dataset
- *N* = 350,000
- *d* = 127

- Webspam dataset
- *N* = 350,000
- *d* = 127



- Webspam dataset
- *N* = 350,000
- *d* = 127





- Webspam dataset
- *N* = 350,000
- *d* = 127








Competitive approximation performance



Competitive approximation performance



Roadmap

- Approximate Bayes review
- Likelihood approximation and dataset compression
- Approximate sufficient statistics
- Accuracy guarantees

Accuracy guarantees

Accuracy guarantees



Accuracy guarantees

• Problem: existing scalable methods lack accuracy guarantees



Question: how do we measure closeness of the exact and approximate posteriors?





Question: how do we measure closeness of the exact and approximate posteriors?





- Question: how do we measure closeness of the exact and approximate posteriors?
- **Recall:** want to compute means, variances, tail probabilities, etc.





- Question: how do we measure closeness of the exact and approximate posteriors?
- **Recall:** want to compute means, variances, tail probabilities, etc.
- Good choice of measure: **1- and 2-Wasserstein distances** d_W





- Question: how do we measure closeness of the exact and approximate posteriors?
- **Recall:** want to compute means, variances, tail probabilities, etc.
- Good choice of measure: **1- and 2-Wasserstein distances** d_W
- Why? $d_W(p, q)$ small implies





- Question: how do we measure closeness of the exact and approximate posteriors?
- **Recall:** want to compute means, variances, tail probabilities, etc.
- Good choice of measure: **1- and 2-Wasserstein distances** d_W
- Why? $d_W(p, q)$ small implies
 - means and variances close





- Question: how do we measure closeness of the exact and approximate posteriors?
- **Recall:** want to compute means, variances, tail probabilities, etc.
- Good choice of measure: **1- and 2-Wasserstein distances** d_W
- Why? $d_W(p, q)$ small implies
 - means and variances close \checkmark
 - (smoothed) tail probabilities close

 $\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z}$

$\pi(\theta Y) =$	$p(Y \theta)\pi_0(\theta)$	$\tilde{\pi}(\boldsymbol{\theta} \mathbf{V}) =$	$e^{\ell(\theta,g(Y))}\pi_0(\theta)$
	\overline{Z}	$\pi(0 1) -$	\tilde{Z}

 $\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z} \qquad \tilde{\pi}(\theta|Y) = \frac{e^{\ell(\theta,g(Y))}\pi_0(\theta)}{\tilde{Z}}$

 $\varepsilon(\theta) = \|\nabla_{\theta} \log p(Y|\theta) - \nabla_{\theta} \ell(\theta, g(Y))\|_{2}$

$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z} \qquad \tilde{\pi}(\theta|Y) = \frac{e^{\ell(\theta,g(Y))}\pi_0(\theta)}{\tilde{Z}}$$

 $\varepsilon(\theta) = \|\nabla_{\theta} \log p(Y|\theta) - \nabla_{\theta} \ell(\theta, g(Y))\|_{2}$

Theorem (H. & Zou 2017, H. 2018). Assume

$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z} \qquad \tilde{\pi}(\theta|Y) = \frac{e^{\ell(\theta,g(Y))}\pi_0(\theta)}{\tilde{Z}}$$

 $\varepsilon(\theta) = \|\nabla_{\theta} \log p(Y|\theta) - \nabla_{\theta} \ell(\theta, g(Y))\|_{2}$

Theorem (H. & Zou 2017, H. 2018). Assume

· π is "well-behaved" and

$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z} \qquad \tilde{\pi}(\theta|Y) = \frac{e^{\ell(\theta,g(Y))}\pi_0(\theta)}{\tilde{Z}}$$

 $\varepsilon(\theta) = \|\nabla_{\theta} \log p(Y|\theta) - \nabla_{\theta} \ell(\theta, g(Y))\|_{2}$

Theorem (H. & Zou 2017, H. 2018). Assume

- · π is "well-behaved" and
- $\varepsilon(\theta) \leq \varepsilon$.

$$\pi(\theta|Y) = \frac{p(Y|\theta)\pi_0(\theta)}{Z} \qquad \tilde{\pi}(\theta|Y) = \frac{e^{\ell(\theta,g(Y))}\pi_0(\theta)}{\tilde{Z}}$$

 $\varepsilon(\theta) = \|\nabla_{\theta} \log p(Y|\theta) - \nabla_{\theta} \ell(\theta, g(Y))\|_{2}$

Theorem (H. & Zou 2017, H. 2018). Assume

- · π is "well-behaved" and
- $\varepsilon(\theta) \leq \varepsilon$.

Then $d_W(\pi, \tilde{\pi}) \leq c_\pi \varepsilon$.

 $\tilde{\pi}_M = \operatorname{order} M \operatorname{PASS-LR} \operatorname{approximate} \operatorname{posterior}$

 $\tilde{\pi}_M = \operatorname{order} M \operatorname{PASS-LR} \operatorname{approximate} \operatorname{posterior}$

Theorem (H., Adams, Broderick 2017).

 $\tilde{\pi}_M = \operatorname{order} M \operatorname{PASS-LR} \operatorname{approximate} \operatorname{posterior}$

Theorem (H., Adams, Broderick 2017). Assume prior and data are "well-behaved".

 $\tilde{\pi}_M = \operatorname{order} M \operatorname{PASS-LR} \operatorname{approximate} \operatorname{posterior}$

Theorem (H., Adams, Broderick 2017). Assume prior and data are "well-behaved".

Then there exist c > 0 and 0 < r < 1 such that

 $\tilde{\pi}_M = \operatorname{order} M \operatorname{PASS-LR} \operatorname{approximate} \operatorname{posterior}$

Theorem (H., Adams, Broderick 2017). Assume prior and data are "well-behaved".

Then there exist c > 0 and 0 < r < 1 such that

 $d_W(\pi, \tilde{\pi}_M) \le c dr^M$

 $\tilde{\pi}_M = \operatorname{order} M \operatorname{PASS-LR} \operatorname{approximate} \operatorname{posterior}$

Theorem (H., Adams, Broderick 2017). Assume prior and data are "well-behaved".

Then there exist c > 0 and 0 < r < 1 such that

$$d_W(\pi, \tilde{\pi}_M) \le c dr^M$$

• Similar results for other GLMs

• Empirical work

- Empirical work
 - Hierarchical models

- Empirical work
 - Hierarchical models
 - Applications to a wider range of GLMs

- Empirical work
 - Hierarchical models
 - Applications to a wider range of GLMs
 - Practitioner buy-in

- Empirical work
 - Hierarchical models
 - Applications to a wider range of GLMs
 - Practitioner buy-in
- Very- and ultra-high dimensional parameter spaces

- Empirical work
 - Hierarchical models
 - Applications to a wider range of GLMs
 - Practitioner buy-in
- Very- and ultra-high dimensional parameter spaces
- Non-parametric models:

- Empirical work
 - Hierarchical models
 - Applications to a wider range of GLMs
 - Practitioner buy-in
- Very- and ultra-high dimensional parameter spaces
- Non-parametric models:
 - Coresets for Gaussian processes, connections to inducing point methods

- Empirical work
 - Hierarchical models
 - Applications to a wider range of GLMs
 - Practitioner buy-in
- Very- and ultra-high dimensional parameter spaces
- Non-parametric models:
 - Coresets for Gaussian processes, connections to inducing point methods
- Combinatorial parameter spaces

Thanks!

J. H. Huggins, R. P. Adams, T. Broderick. PASS-GLM: polynomial approximate sufficient statistics for scalable Bayesian GLM inference. NIPS, 2017.

J. H. Huggins*, J. Zou*. Quantifying the Accuracy of Approximate Diffusions and Markov Chains. AISTATS, 2017.

J. H. Huggins, T. Campbell, T. Broderick. Coresets for Scalable Bayesian Logistic Regression. NIPS, 2016.




$$b(\theta) = \nabla_{\theta} \log \pi(\theta \,|\, Y)$$
$$\tilde{b}(\theta) = \nabla_{\theta} \log \tilde{\pi}(\theta \,|\, Y)$$



$$b(\theta) = \nabla_{\theta} \log \pi(\theta \mid Y)$$
$$\tilde{b}(\theta) = \nabla_{\theta} \log \tilde{\pi}(\theta \mid Y)$$

$$d\theta_t = b(\theta_t)dt + \sqrt{2}dW_t$$
$$d\tilde{\theta}_t = \tilde{b}(\tilde{\theta}_t)dt + \sqrt{2}d\tilde{W}_t$$

$$b(\theta) = \nabla_{\theta} \log \pi(\theta \mid Y)$$
$$\tilde{b}(\theta) = \nabla_{\theta} \log \tilde{\pi}(\theta \mid Y)$$

$$d\theta_t = b(\theta_t)dt + \sqrt{2}dW_t$$
$$d\tilde{\theta}_t = \tilde{b}(\tilde{\theta}_t)dt + \sqrt{2}d\tilde{W}_t$$

$$b(\theta) = \nabla_{\theta} \log \pi(\theta \mid Y)$$
$$\tilde{b}(\theta) = \nabla_{\theta} \log \tilde{\pi}(\theta \mid Y)$$

$$d\theta_t = b(\theta_t)dt + \sqrt{2}dW_t$$
$$d\tilde{\theta}_t = \tilde{b}(\tilde{\theta}_t)dt + \sqrt{2}dW_t$$

$$b(\theta) = \nabla_{\theta} \log \pi(\theta \mid Y)$$
$$\tilde{b}(\theta) = \nabla_{\theta} \log \tilde{\pi}(\theta \mid Y)$$

$$d\theta_t = b(\theta_t)dt + \sqrt{2}dW_t$$
$$d\theta_t = \tilde{b}(\theta_t)dt + \sqrt{2}dW_t$$
$$d\theta_t = \tilde{b}(\theta_t)dt + \sqrt{2}dW_t$$

$$b(\theta) = \nabla_{\theta} \log \pi(\theta \mid Y)$$
$$\tilde{b}(\theta) = \nabla_{\theta} \log \tilde{\pi}(\theta \mid Y)$$

$$d\theta_{t} = b(\theta_{t})dt + \sqrt{2}dW_{t}$$
$$d\theta_{t} = \tilde{b}(\theta_{t})dt + \sqrt{2}dW_{t}$$
$$d\theta_{t} = \tilde{b}(\theta_{t})dt + \sqrt{2}dW_{t}$$

$$b(\theta) = \nabla_{\theta} \log \pi(\theta \mid Y)$$
$$\tilde{b}(\theta) = \nabla_{\theta} \log \tilde{\pi}(\theta \mid Y)$$

$$d\theta_{t} = b(\theta_{t})dt + \sqrt{2}dW_{t}$$
$$d\theta_{t} = \tilde{b}(\theta_{t})dt + \sqrt{2}dW_{t}$$
$$d\theta_{t} = \tilde{b}(\theta_{t})dt + \sqrt{2}dW_{t}$$

$$b(\theta) = \nabla_{\theta} \log \pi(\theta \mid Y)$$
$$\tilde{b}(\theta) = \nabla_{\theta} \log \tilde{\pi}(\theta \mid Y)$$

$$\begin{split} \tilde{b}(\theta_t) dt & \sqrt{2} dW_t \\ \tilde{b}(\theta_t) dt & \theta_{t+dt} \\ \theta_t & \theta_{t+dt} \\ d\theta_t & \theta_{t+dt} \\ d\theta_t & = b(\theta_t) dt + \sqrt{2} dW_t \\ d\tilde{\theta}_t &= \tilde{b}(\tilde{\theta}_t) dt + \sqrt{2} d\tilde{W}_t \end{split}$$

- Diffusion: continuous-time
 Markov process with unique
 stationary distribution
- Intuition: if diffusion mixes quickly, then gradient errors don't have time to build up

$$b(\theta) = \nabla_{\theta} \log \pi(\theta \,|\, Y)$$
$$\tilde{b}(\theta) = \nabla_{\theta} \log \tilde{\pi}(\theta \,|\, Y)$$



Theorem (H. & Zou 2017, H. 2018).

Theorem (H. & Zou 2017, H. 2018).

Assume the diffusion converges at rate r(t).



Theorem (H. & Zou 2017, H. 2018).

Assume the diffusion converges at rate r(t).

Let $I(r) = \int r(t) dt$.



Theorem (H. & Zou 2017, H. 2018).

Assume the diffusion converges at rate r(t).

Let $I(r) = \int r(t) dt$.

Then $d_W(\pi, \tilde{\pi}) \leq I(r) \mathbb{E}_{\tilde{\pi}}[\|b - \tilde{b}\|_2].$



Theorem (H. & Zou 2017, H. 2018).

Assume the diffusion converges at rate r(t).

Let $I(r) = \int r(t) dt$.

Then $d_W(\pi, \tilde{\pi}) \leq \underbrace{I(r)}_{c_{\pi}} \mathbb{E}_{\tilde{\pi}}[\|b - \tilde{b}\|_2].$



Theorem (H. & Zou 2017, H. 2018).

Assume the diffusion converges at rate r(t).

Let $I(r) = \int r(t) dt$.

Then $d_W(\pi, \tilde{\pi}) \leq \underbrace{I(r)}_{c_{\pi}} \underbrace{\mathbb{E}_{\tilde{\pi}}[\|b - \tilde{b}\|_2]}_{\varepsilon}$.



Theorem (H. & Zou 2017, H. 2018).

Assume the diffusion converges at rate r(t).

Let $I(r) = \int r(t) dt$.

Then
$$d_W(\pi, \tilde{\pi}) \leq \underbrace{I(r)}_{c_{\pi}} \underbrace{\mathbb{E}_{\tilde{\pi}}[\|b - \tilde{b}\|_2]}_{\varepsilon}$$
.

Proof techniques:



Theorem (H. & Zou 2017, H. 2018).

Assume the diffusion converges at rate r(t).

Let $I(r) = \int r(t) dt$.

Then
$$d_W(\pi, \tilde{\pi}) \leq \underbrace{I(r)}_{c_{\pi}} \underbrace{\mathbb{E}_{\tilde{\pi}}[\|b - \tilde{b}\|_2]}_{\varepsilon}$$
.

- Proof techniques:
 - Stein's method (for 1-Wasserstein version)



Theorem (H. & Zou 2017, H. 2018).

Assume the diffusion converges at rate r(t).

Let $I(r) = \int r(t) dt$.

Then
$$d_W(\pi, \tilde{\pi}) \leq \underbrace{I(r)}_{c_{\pi}} \underbrace{\mathbb{E}_{\tilde{\pi}}[\|b - \tilde{b}\|_2]}_{\varepsilon}$$
.

 $i = \frac{1}{1} \int_{1}^{2} \frac{r(t)}{1} \int_{10}^{1} \frac{r(t)}{10} \int_{20}^{1} \frac{r(t)}{10} \int_{10}^{1} \frac{r(t)}{10} \int_{10}^{1}$

- Proof techniques:
 - Stein's method (for 1-Wasserstein version)
 - A coupling argument + Ito's lemma (for 2-Wasserstein version)