



Motivation & Contribution

- 1) Improving stationary mixing efficiency in SGMCMC by leveraging a generalized (potentially heavy-tailed) kinetics.
- 2) Alleviating numerical issue and satisfying conditions for stationarity by leveraging smooth version of generalized kinetics.
- 3) Ameliorating convergence issue by introducing additional first order dynamics and stochastic resampling.

MGHMC

Hamiltonian Monte Carlo (HMC) leverages Hamiltonian dynamics to propose new samples for x from $p(x) \propto \exp(-U(x))$, driven by the following partial differential equations (PDE):

$$\frac{dx}{dt} = \nabla_p K(p) \quad , \quad \frac{dp}{dt} = -\nabla_x U(x) \quad .$$

- Consider the generalized kinetic $K(p; m, a) = \frac{|p|^{1/a}}{m}$, $a, m > 0$
- monomial Gamma (MG) distribution:

$$\pi(p; m, a) = \exp(-K(p; m, a)) = \frac{m^{-a}}{2\Gamma(a+1)} e^{-\frac{|p|^{1/a}}{m}} \quad .$$

Theorem

For univariate target distribution, the one time lag autocorrelation $\rho(x_t, x_{t+1})$ of the analytic MG-SS parameterized by a asymptotically approaches zero when $a \rightarrow \infty$, under regularity condition of $U(x)$ and stationary assumption.

- In addition to above, the MG-HMC with large a is particularly advantageous for sampling multimodal distributions.
- Such a performance gain does not come in free.

SGMCMC

- Sampling from $f(\theta) \propto \exp(-U(\theta))$ using minibatch data.
- SGHMC (stochastic gradient Hamiltonian Monte Carlo)

$$d\theta = \partial_p K(p) dt$$

$$dp = -\partial_\theta \tilde{U}(\theta) dt - B(\theta) \partial_p K(p) dt + \mathcal{N}(0, 2B(\theta) dt) \quad .$$

- SGNHT (stochastic gradient Nosé-Hoover thermostat)

$$d\theta = \partial_p K(p) dt$$

$$dp = -\partial_\theta \tilde{U}(\theta) dt - \xi \partial_p K(p) dt + \mathcal{N}(0, 2B(\theta) dt)$$

$$d\xi = (p^T p - 1) dt \quad .$$

SGMGS (Stochastic gradient MG sampler)

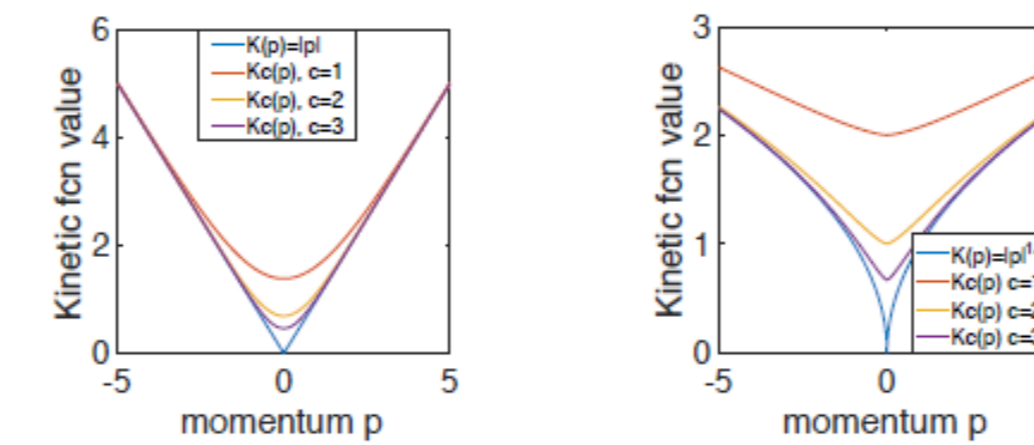
- Consider applying generalized kinetics $K(p)$ to SGNHT.

$$d\theta = \nabla K(p) dt \quad ,$$

$$dp = -[\nabla \tilde{U}(\theta) + \xi \odot \nabla K(p)] dt + \sqrt{2A} dW \quad ,$$

$$d\xi = (\nabla K(p) \odot \nabla K(p) - \nabla^2 K(p)) dt \quad .$$

- The existence and uniqueness of the solutions to the Fokker-Planck equation require Lipschitz continuity of drift and diffusion vectors.
- We propose a softened kinetics for $a = \{1, 2\}$



Convergence issue and remedies

- Consider a Hamiltonian system defined in a more general form

$$H = K(p) + U(\theta) + F(\xi) \quad ,$$

- We consider adding Brownian motion to θ and ξ

$$d\theta = \nabla K_c(p) dt - \sigma_\theta \nabla U(\theta) dt + \sqrt{2\sigma_\theta} dW$$

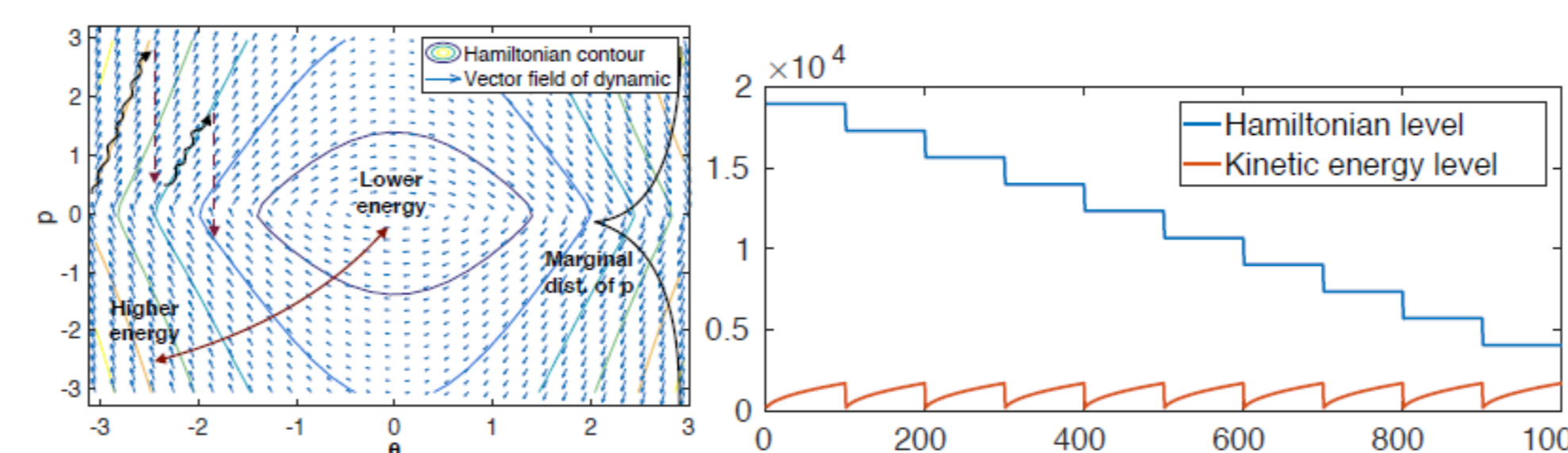
$$dp = -(\sigma_p + \gamma \nabla F(\xi)) \odot \nabla K_c(p) dt$$

$$- \nabla U(\theta) dt + \sqrt{2\sigma_p} dW \quad ,$$

$$d\xi = \gamma [\nabla K_c(p) \odot \nabla K_c(p) - \nabla^2 K_c(p)] dt$$

$$- \sigma_\xi \nabla F(\xi) dt + \sqrt{2\sigma_\xi} dW \quad .$$

- $-\sigma_\theta \nabla U(\theta) dt + \sqrt{2\sigma_\theta} dW$, compensate for the weak updating signal from $\nabla K(p) = \frac{|p|^{1/a-1}}{am}$, by an immediate gradient
- The proposed SDE has better theoretic guarantee on the existence and convergence of bounded solutions for a
- Resampling makes the sampler to immediately move to a lower Hamiltonian energy level.



Theorem

$$\text{Bias: } \left| \mathbb{E} \hat{\phi}_T - \bar{\phi} \right| = O(1/(Th) + h) \quad ,$$

$$\text{MSE: } \mathbb{E} (\hat{\phi} - \bar{\phi})^2 = O(1/(Th) + h^2) \quad .$$

Multiple-well Synthetic Potential

- Generate samples from a complex multimodal distribution.

$$U(\theta) \triangleq e^{\frac{3}{4}\theta^2 - \frac{3}{2}\sum_{i=1}^{10} c_i \sin(\frac{1}{4}\pi i(\theta+4))}$$

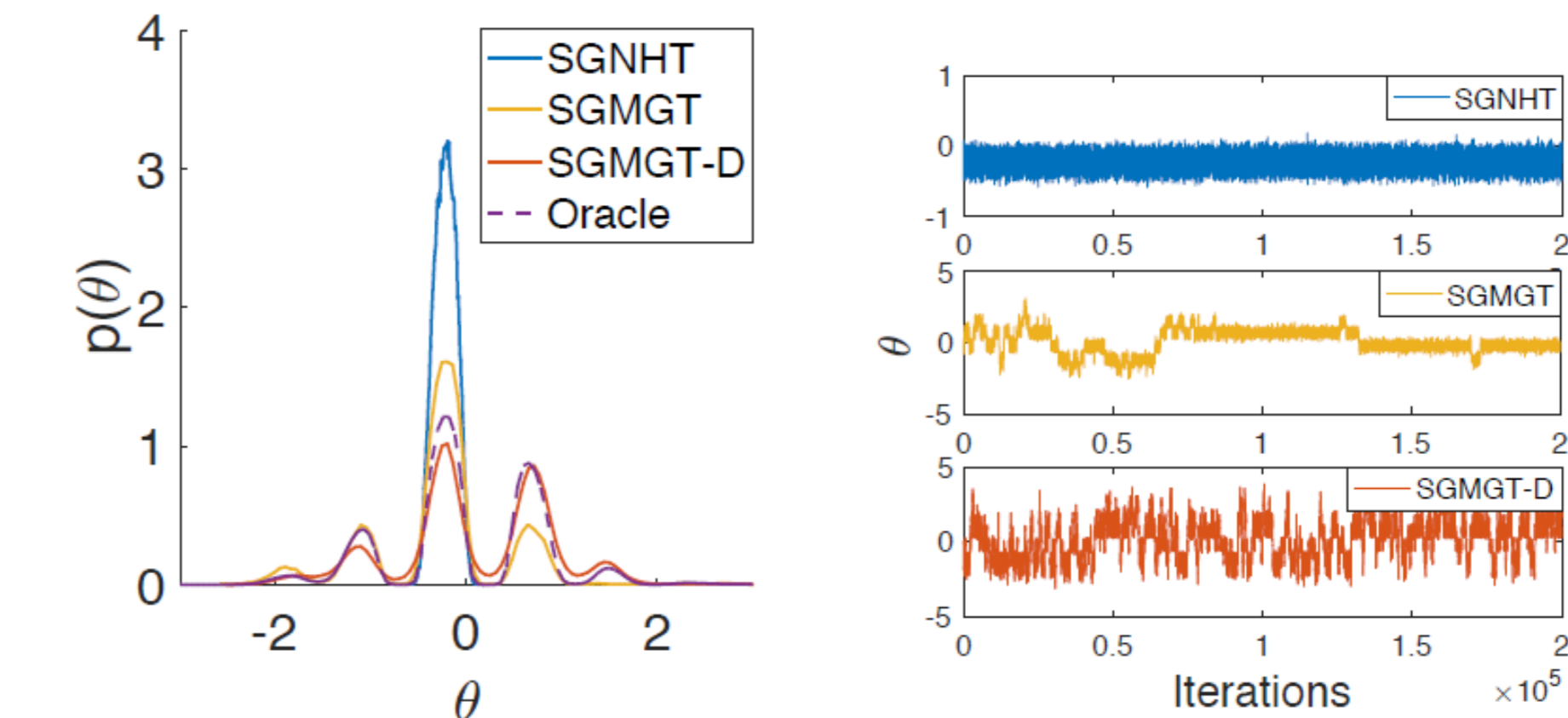


Figure: Synthetic multimodal distribution. Left: empirical distributions for different methods. Right: traceplot for each method.

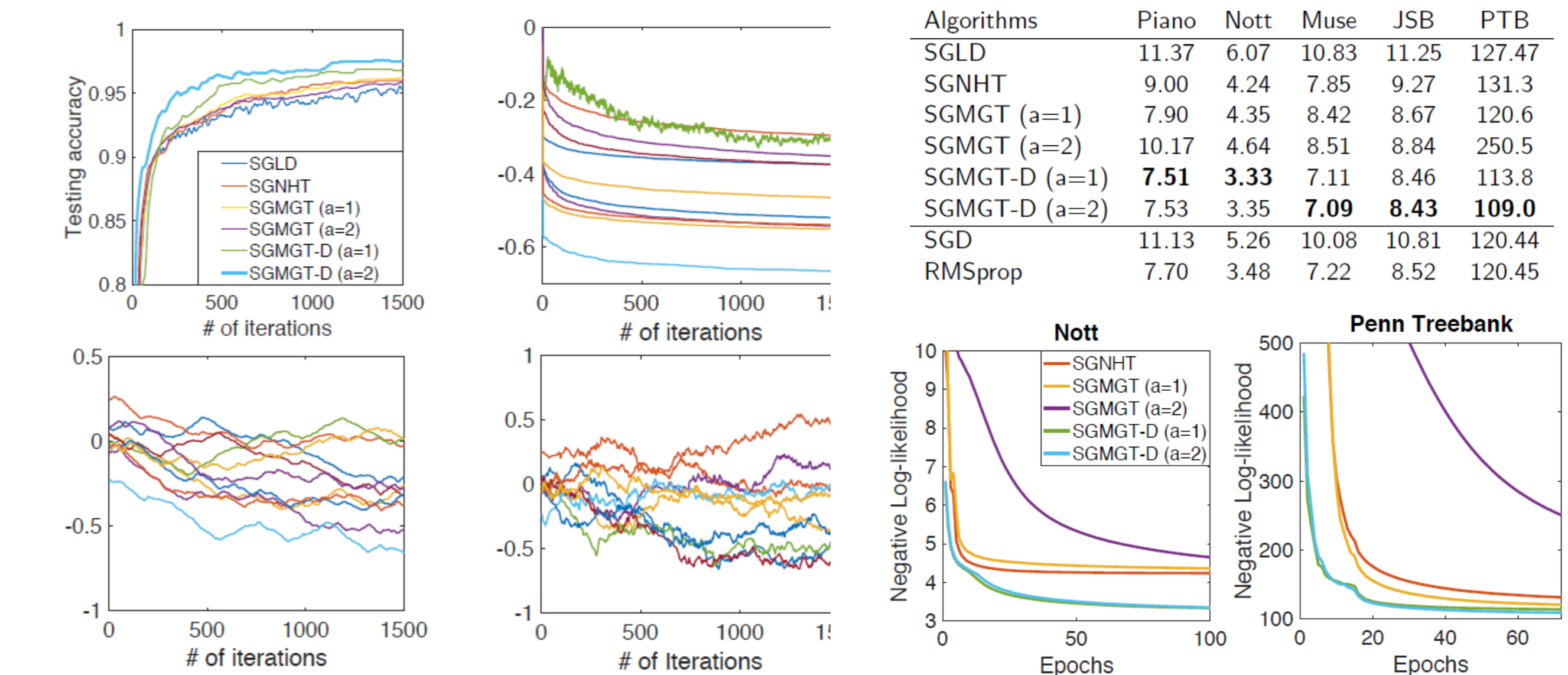
Bayesian Logistic Regression

Table: Average AUROC and median ESS. Dataset dimensionality is indicated in parenthesis after the name of each dataset.

AUROC (D)	A (15)	G (25)	H (14)	P (8)	R (7)	C (87)
SGNHT	0.89	0.75	0.90	0.86	0.95	0.65
SGMGT(a=1)	0.92	0.78	0.91	0.86	0.87	0.70
SGMGT-D(a=1)	0.95	0.86	0.95	0.93	0.98	0.73
SGMGT(a=2)	0.93	0.79	0.93	0.88	0.86	0.62
SGMGT-D(a=2)	0.95	0.90	0.95	0.90	0.97	0.69
ESS (D)	A (15)	G (25)	H (14)	P (8)	R (7)	C (87)
SGNHT	869	941	1911	2077	1761	1873
SGMGT-D(a=1)	3147	2131	2448	4244	1494	3605
SGMGT-D(a=2)	2700	1989	2768	3430	2265	2969

Deep models

- Discriminative RBM & Recurrent Neural Networks. Assuming flat prior. 5000 MC samples.



Conclusion

- Scalable MCMC inference with generalized HMC variants.

Future direction:

- Adaptive selection of monomial parameters
- Connection to optimization methods.