

An energy-based model approach to the estimation of rare event probabilities

Lea Friedli

Joint work with:

Arnaud Doucet, University of Oxford

David Ginsbourger, University of Bern

Niklas Linde, University of Lausanne

No risk, no fun
Rare event probabilities

No risk, no fun

Rare event probabilities

- Bayesian inversion inferring property field θ given measurements \mathbf{y}
- Interested in quantity depending on field through $\theta \mapsto \mathcal{R}(\theta)$
 - $\mathbb{P}(\mathcal{R}(\theta) \geq T \mid \mathbf{y})$, risk of failure of a system

No risk, no fun

Rare event probabilities

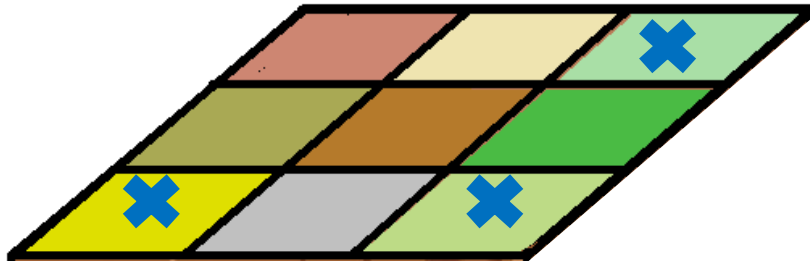
- Bayesian inversion inferring property field θ given measurements \mathbf{y}
- Interested in quantity depending on field through $\theta \mapsto \mathcal{R}(\theta)$
→ $\mathbb{P}(\mathcal{R}(\theta) \geq T \mid \mathbf{y})$, risk of failure of a system
- High-dimensional target space, low-dimensional risk space
- Applications in finance, engineering, environmental science,...

No risk, no fun

Rare event probabilities

- Bayesian inversion inferring property field θ given measurements \mathbf{y}
- Interested in quantity depending on field through $\theta \mapsto \mathcal{R}(\theta)$
→ $\mathbb{P}(\mathcal{R}(\theta) \geq T \mid \mathbf{y})$, risk of failure of a system
- High-dimensional target space, low-dimensional risk space
- Applications in finance, engineering, environmental science,...
- Traditional Monte Carlo approach: excessive number of samples needed
- Variational approach (Valsson & Parrinello 2014) → Energy based models

Analytical toy example

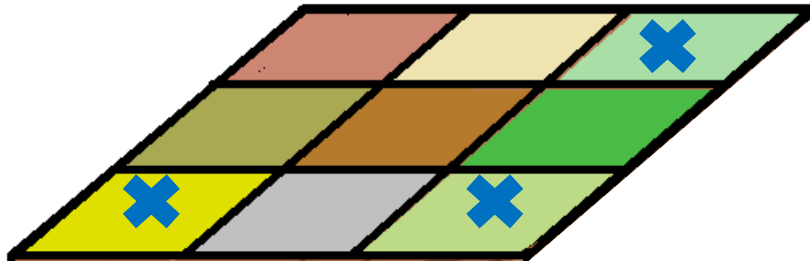


θ : Property field (contamination values)

$\mathbf{y} = \mathcal{G}(\theta) + \epsilon$: (Local) measurements

→ Posterior PDF $p(\theta|\mathbf{y})$

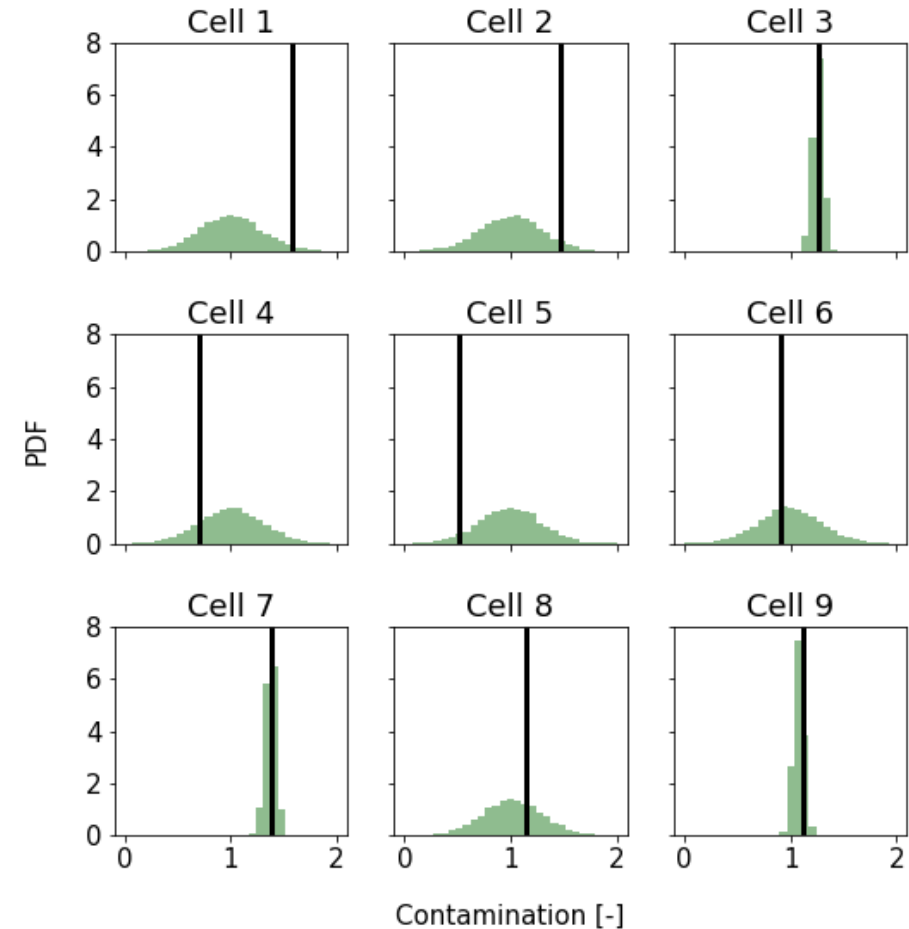
Analytical toy example



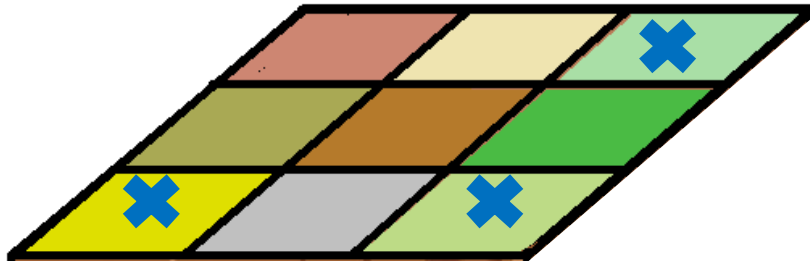
θ : Property field (contamination values)

$y = \mathcal{G}(\theta) + \epsilon$: (Local) measurements

→ Posterior PDF $p(\theta|\mathbf{y})$



Analytical toy example



$\mathcal{R}(\boldsymbol{\theta})$: Quantity of property field

$\rightarrow \mathbb{P}(\mathcal{R}(\boldsymbol{\theta}) \geq T \mid \mathbf{y}) = ?$

$\boldsymbol{\theta}$: Property field (contamination values)

$\mathbf{y} = \mathcal{G}(\boldsymbol{\theta}) + \epsilon$: (Local) measurements

\rightarrow Posterior PDF $p(\boldsymbol{\theta} \mid \mathbf{y})$

Analytical toy example



$\boldsymbol{\theta}$: Property field (contamination values)

$\mathbf{y} = \mathcal{G}(\boldsymbol{\theta}) + \epsilon$: (Local) measurements

→ Posterior PDF $p(\boldsymbol{\theta}|\mathbf{y})$

$\mathcal{R}(\boldsymbol{\theta})$: Quantity of property field

→ $\mathbb{P}(\mathcal{R}(\boldsymbol{\theta}) \geq T | \mathbf{y}) = ?$

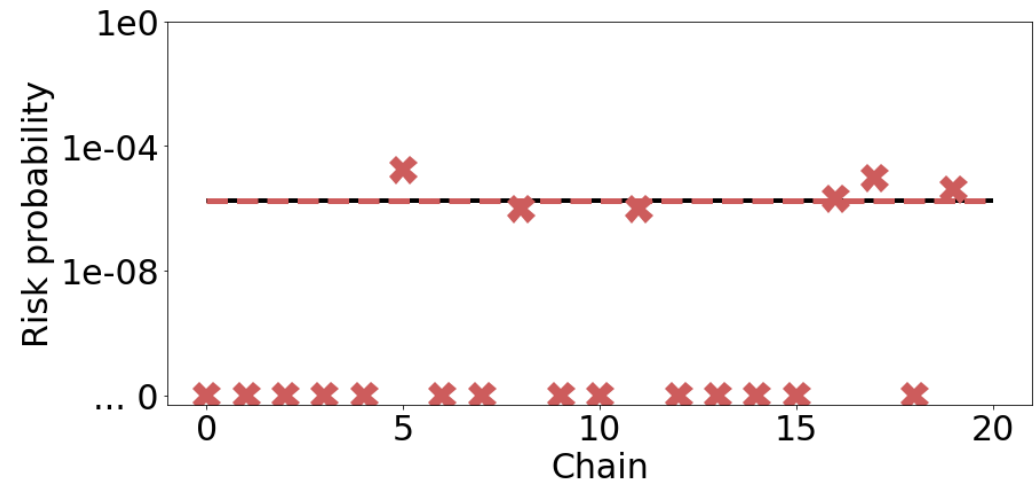
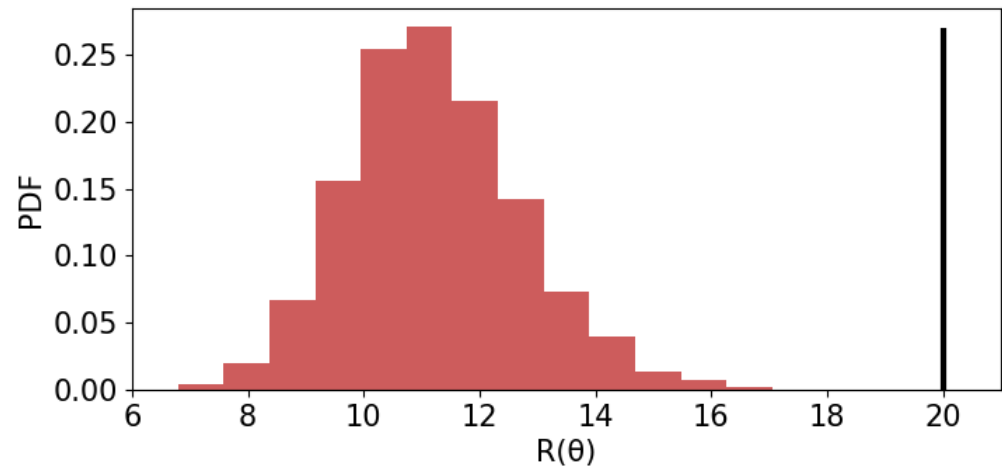
Linear quadratic dose response

$$\mathcal{R}(\boldsymbol{\theta}) = \sum_{m=1}^9 \theta_m^2$$

→ $\mathbb{P}(\mathcal{R}(\boldsymbol{\theta}) \geq 20.0 | \mathbf{y}) = 1.76 \times 10^{-6}$

Metropolis Hastings

- 20 chains, 1 Million iterations each
- Gaussian proposals, AR 30%
- Convergence after 2'000 iterations (Gelman-Rubin)



Energy based model approach

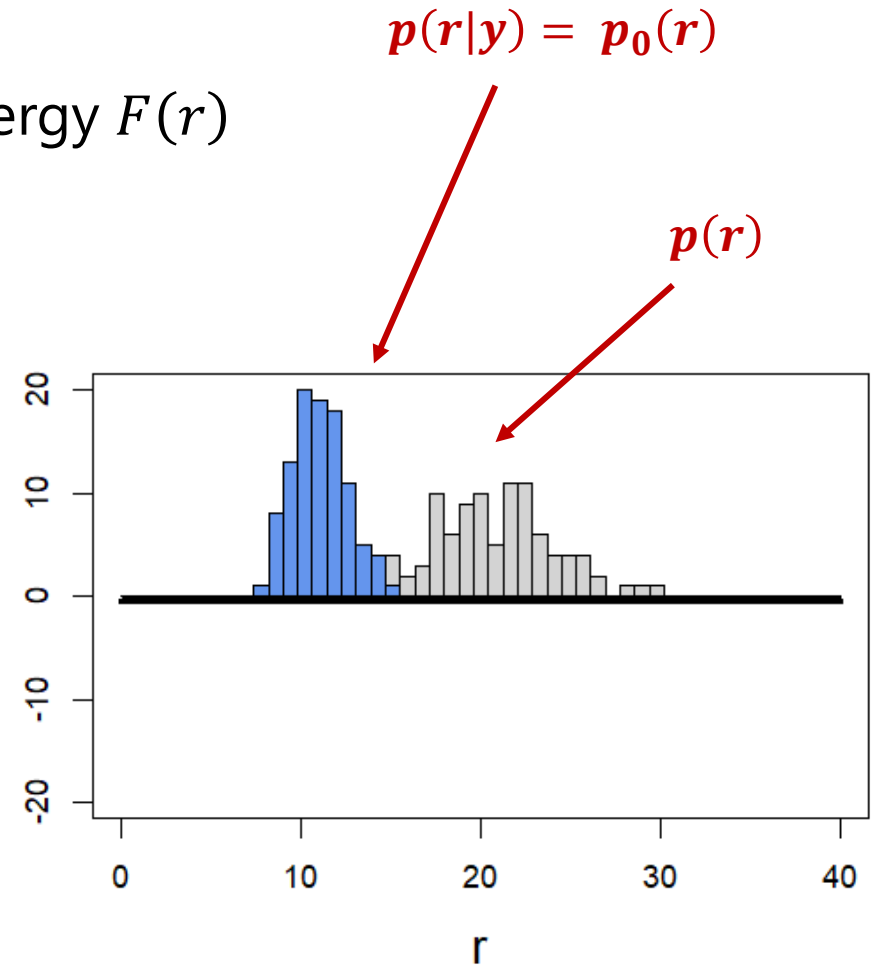
- $\mathbb{P}(\mathcal{R}(\theta) \geq T | y) = \int_T^\infty p(r|y) dr$
- $R = \mathcal{R}(\theta), p(r|y) = \frac{\exp(-F(r))}{\int \exp(-F(s)) ds} \rightarrow$ Find free energy $F(r)$

Energy based model approach

- $\mathbb{P}(\mathcal{R}(\theta) \geq T | y) = \int_T^\infty p(r|y) dr$
- $R = \mathcal{R}(\theta)$, $p(r|y) = \frac{\exp(-F(r))}{\int \exp(-F(s)) ds} \rightarrow$ Find free energy $F(r)$
- Introduce bias potential $V: \mathbb{R} \rightarrow \mathbb{R}, r \mapsto V(r)$

$$p_V(r) = \frac{\exp(-(F(r)+V(r)))}{\int \exp(-(F(s)+V(s))) ds}$$

- Select PDF $p(r)$ with mass on $[T, \infty)$

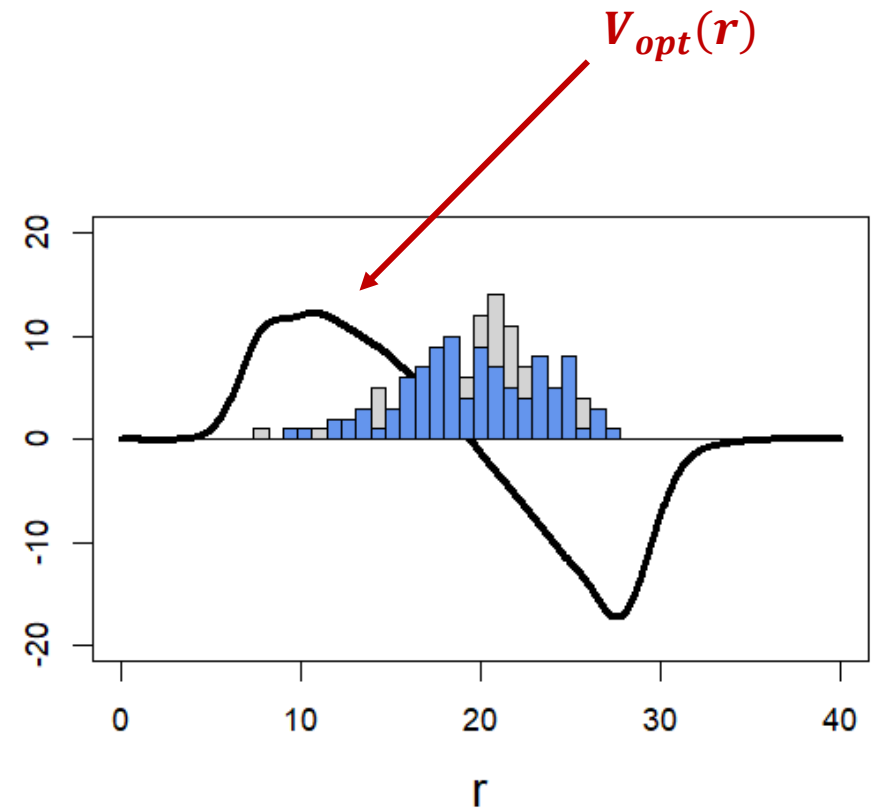


Energy based model approach

- $\mathbb{P}(\mathcal{R}(\theta) \geq T | y) = \int_T^\infty p(r|y) dr$
- $R = \mathcal{R}(\theta)$, $p(r|y) = \frac{\exp(-F(r))}{\int \exp(-F(s)) ds} \rightarrow$ Find free energy $F(r)$
- Introduce bias potential $V: \mathbb{R} \rightarrow \mathbb{R}, r \mapsto V(r)$

$$p_V(r) = \frac{\exp(-(F(r)+V(r)))}{\int \exp(-(F(s)+V(s))) ds}$$

- Select PDF $p(r)$ with mass on $[T, \infty)$
- Find $V_{opt}(r)$ minimizing $V \mapsto KL(p||p_V)$



Energy based model approach

- $\mathbb{P}(\mathcal{R}(\theta) \geq T | y) = \int_T^\infty p(r|y) dr$
- $R = \mathcal{R}(\theta), p(r|y) = \frac{\exp(-F(r))}{\int \exp(-F(s)) ds} \rightarrow$ Find free energy $F(r)$

- Introduce bias potential $V: \mathbb{R} \rightarrow \mathbb{R}, r \mapsto V(r)$

$$p_V(r) = \frac{\exp(-(F(r)+V(r)))}{\int \exp(-(F(s)+V(s))) ds}$$

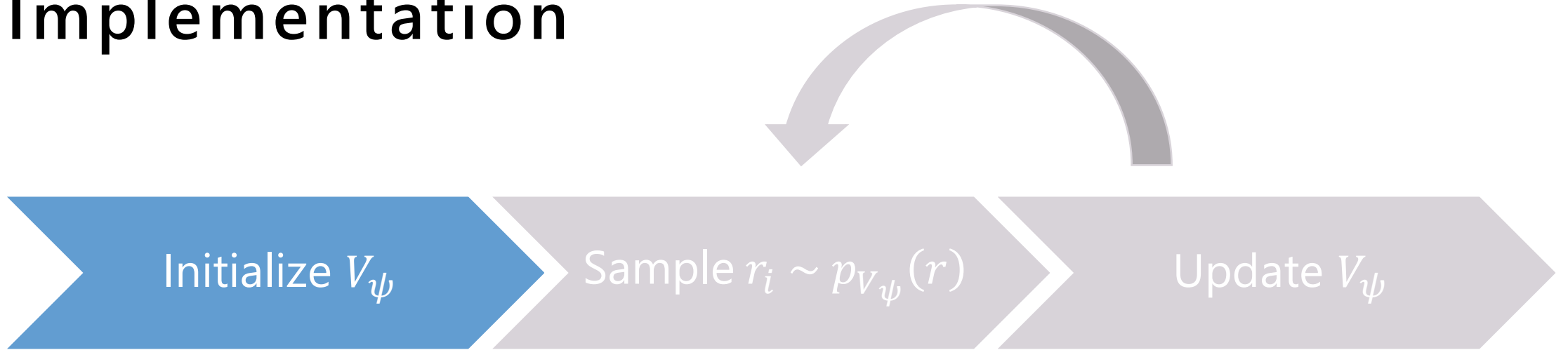
- Select PDF $p(r)$ with mass on $[T, \infty)$
- Find $V_{opt}(r)$ minimizing $V \mapsto KL(p||p_V)$
- $F(r) = -\log(p(r)) - V_{opt}(r)$

Variational approach

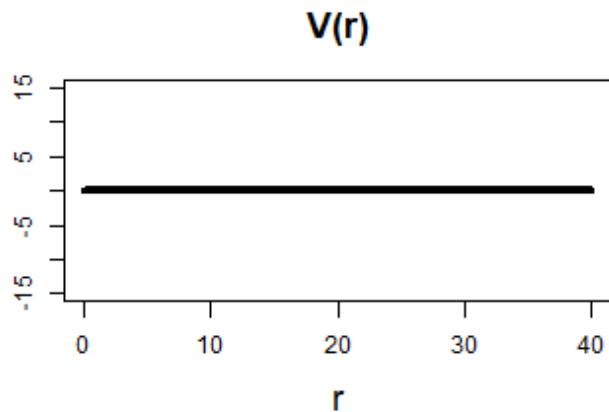
Valsson and Parrinello (2014)

- Loss function
- Rare event probabilities

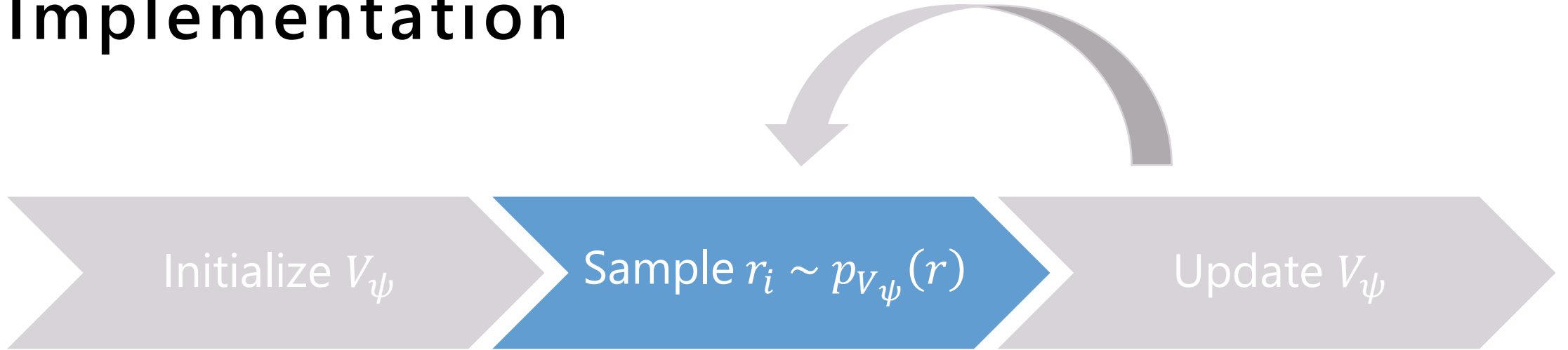
Implementation



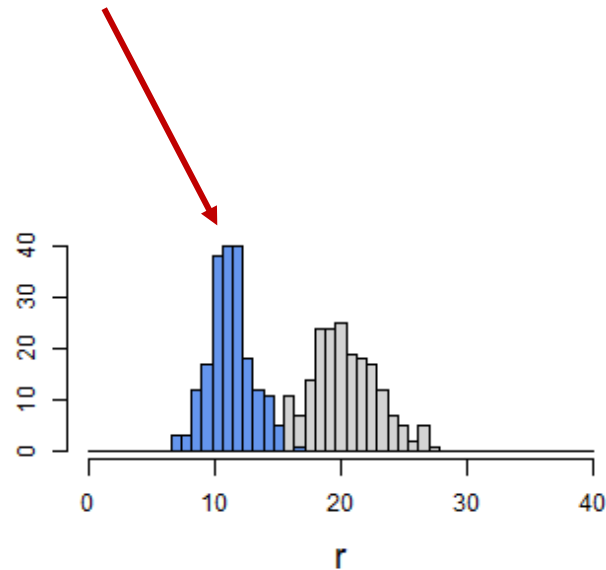
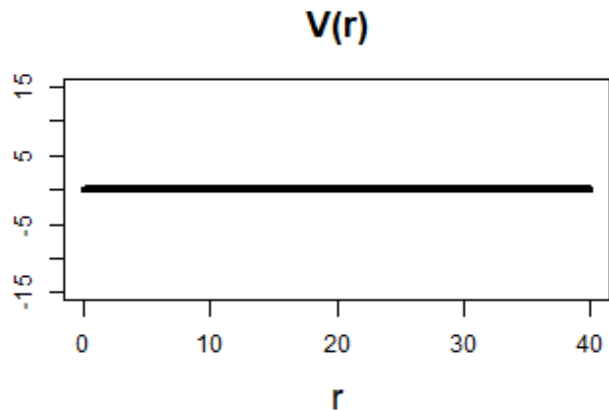
Parameterize bias potential $r \mapsto V_\psi(r)$ with neural network, radial basis functions, splines,...



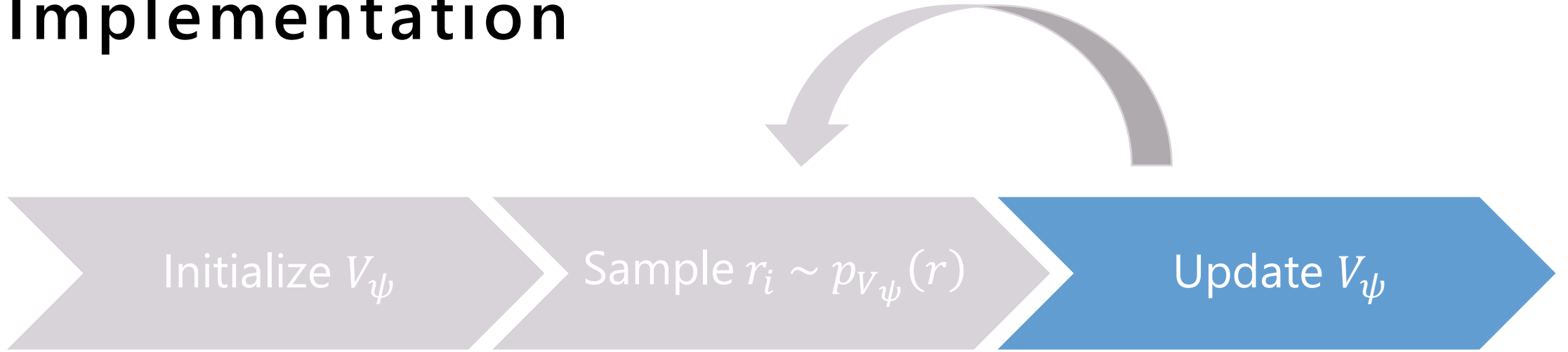
Implementation



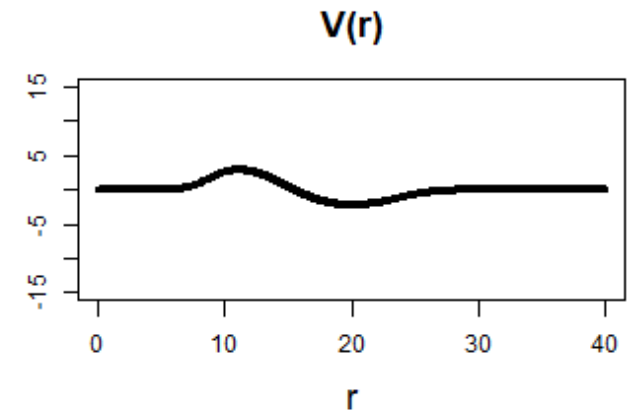
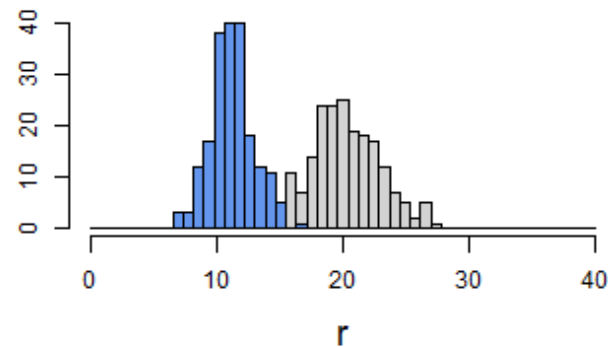
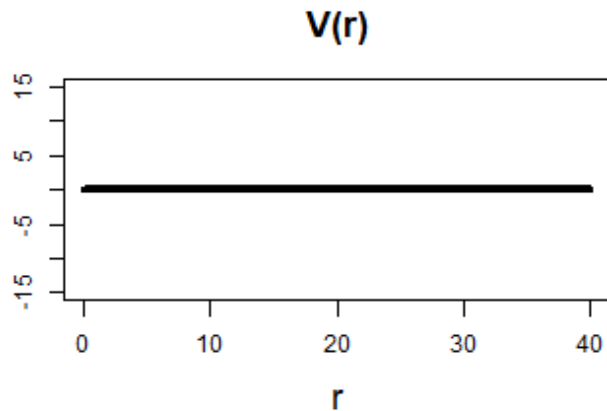
Sample $r_i \sim p_{V_\psi}(r)$ using MCMC



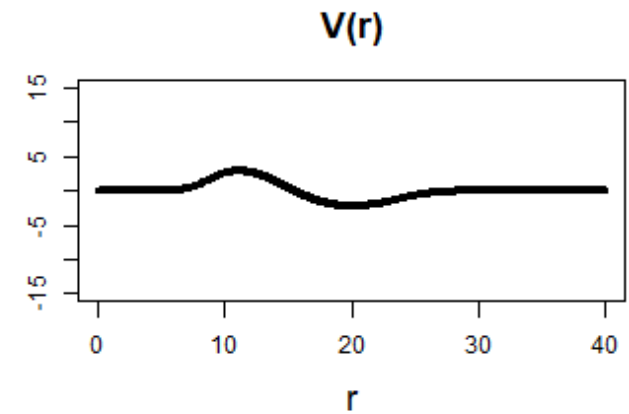
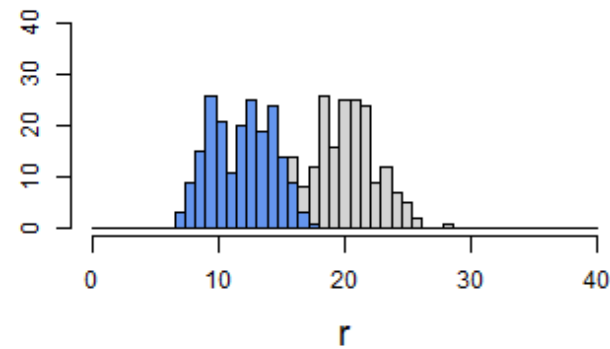
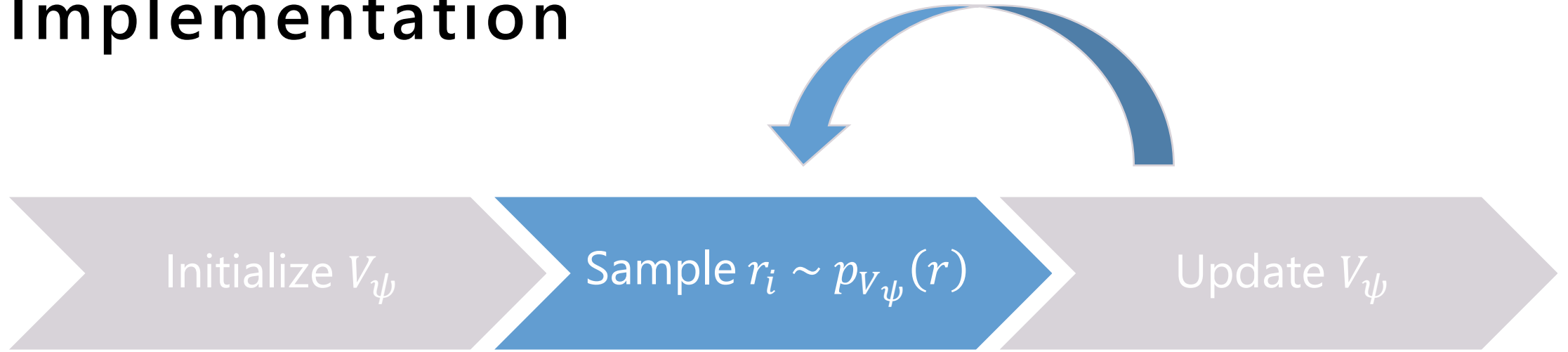
Implementation



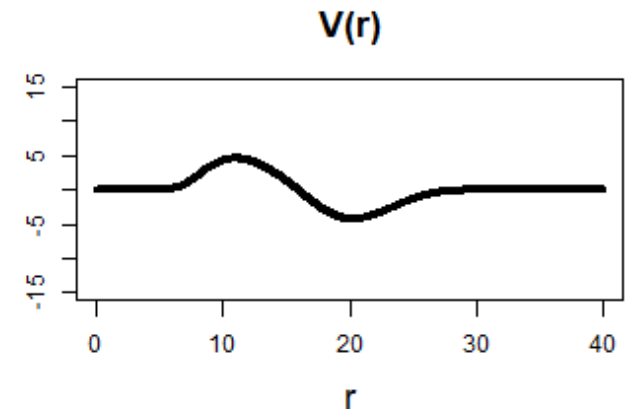
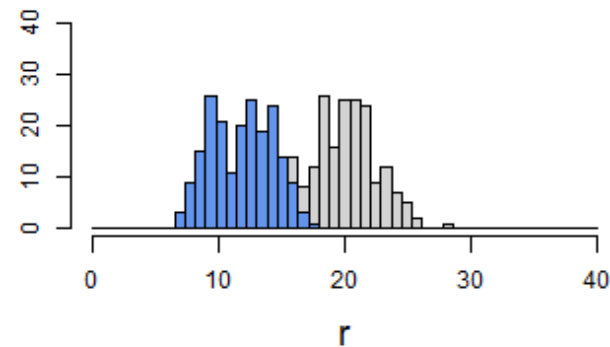
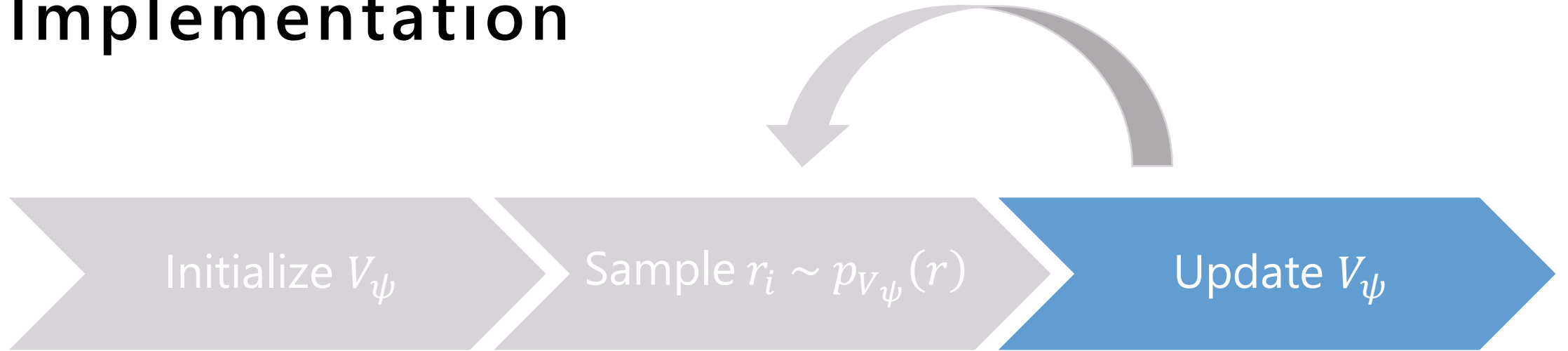
Update ψ using stochastic gradient descent to minimize $KL(p||p_V)$



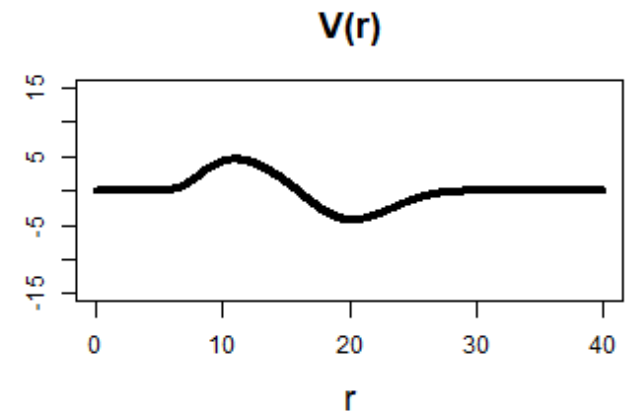
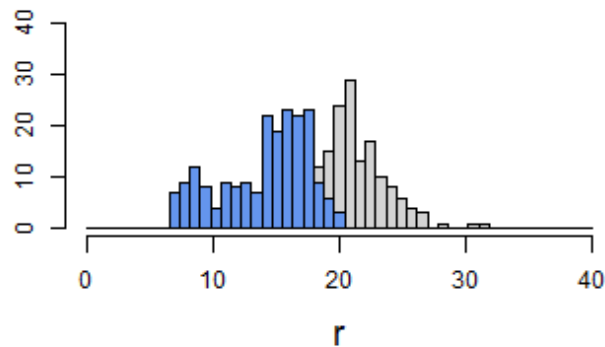
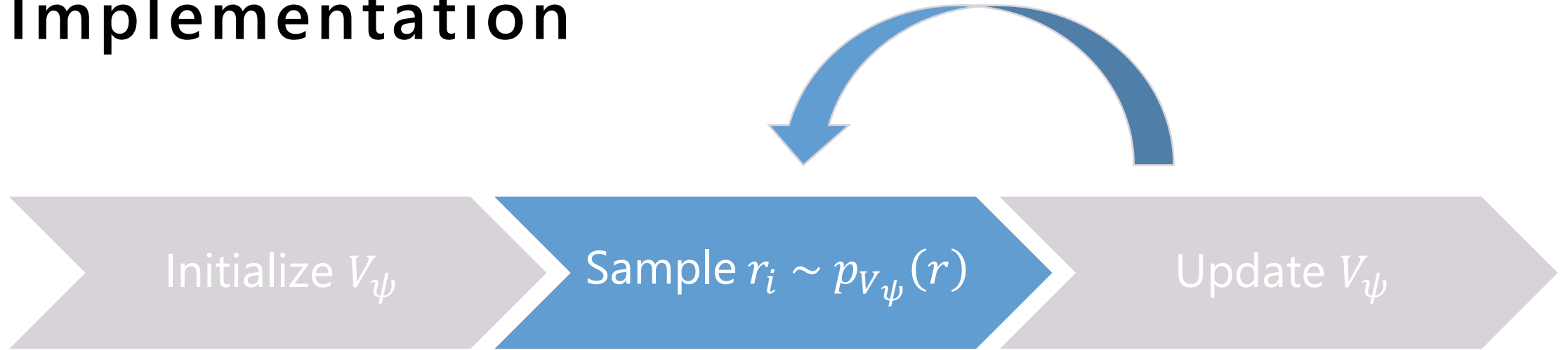
Implementation



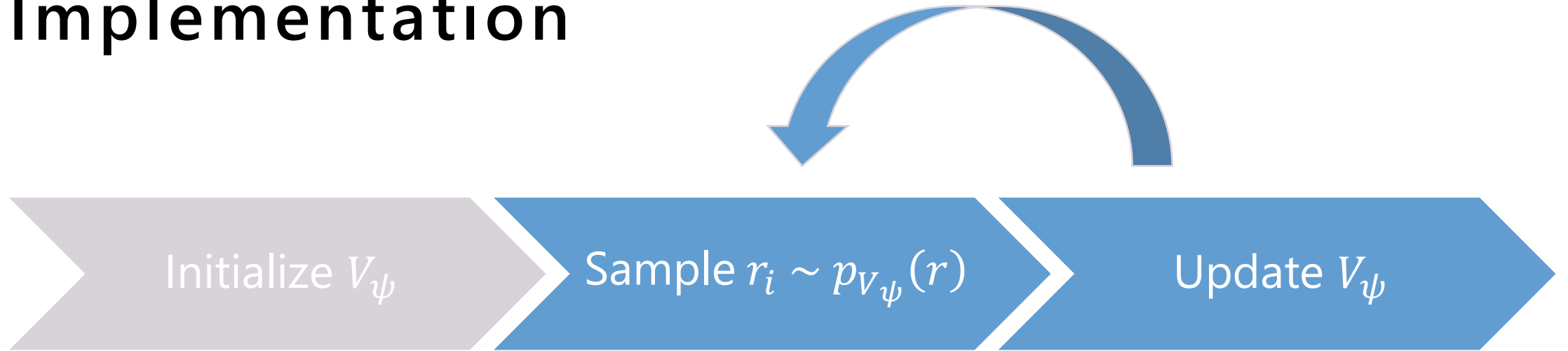
Implementation



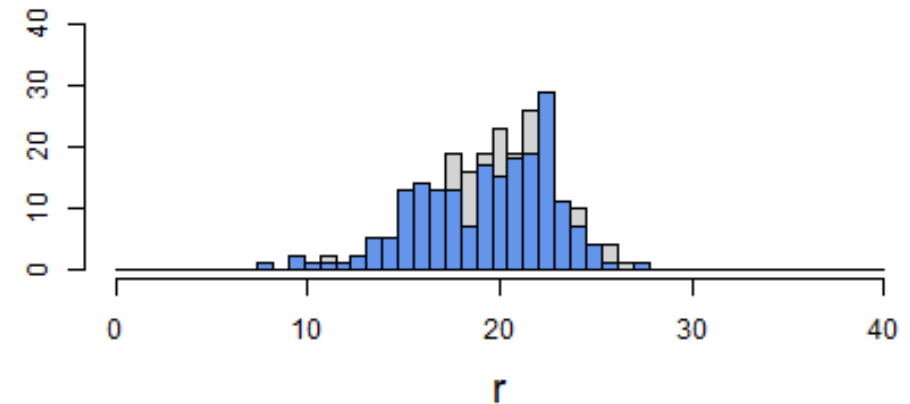
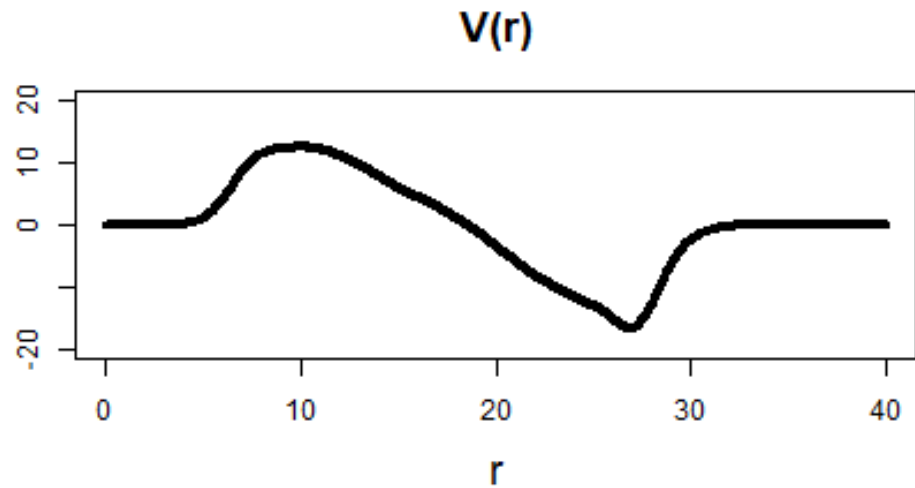
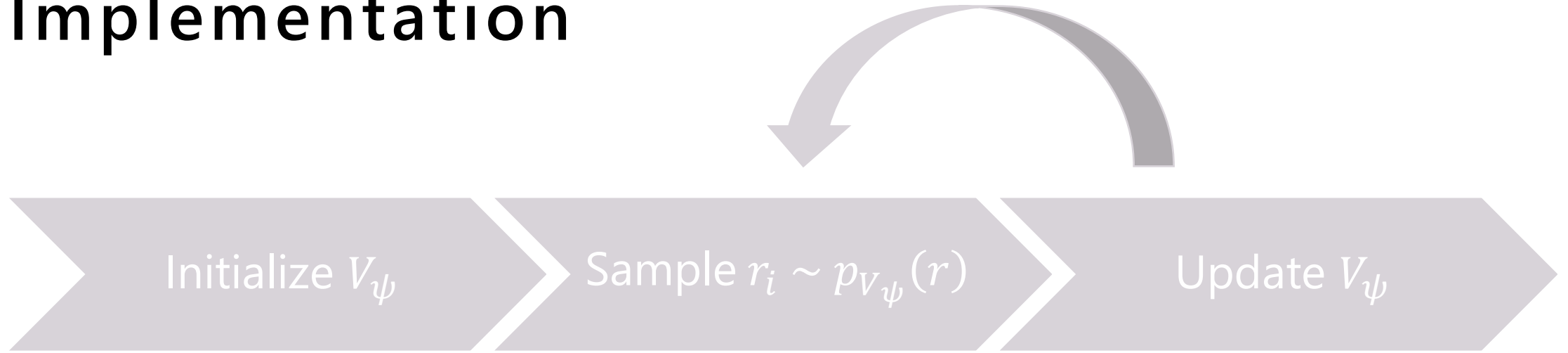
Implementation



Implementation

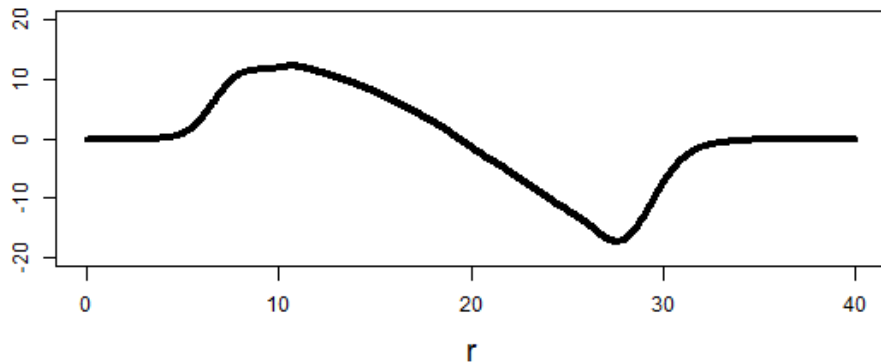


Implementation

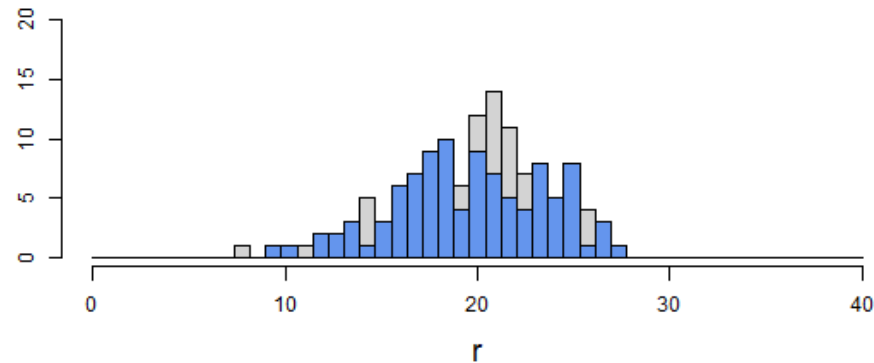


When do we stop?

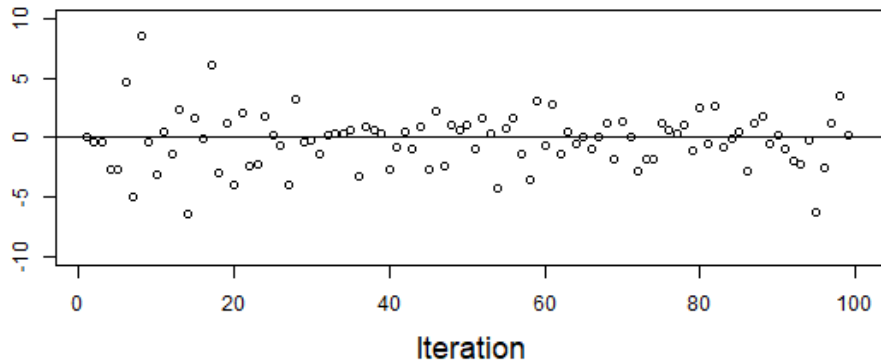
$V(r)$



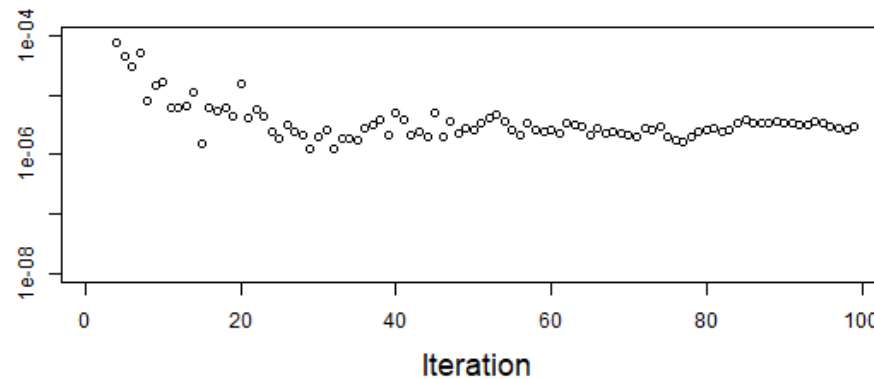
Samples $p(r)$ and $pV(r)$



Loss



Risk probability



Kernel Stein discrepancy

Kernel Stein discrepancy between distributions (Riabiz et al. 2022)

$$KSD(p | p_{V_\psi}) = \sqrt{\frac{1}{n^2} \sum_{i,j=1}^n k_p(x_i, x_j)}, \quad x_i \sim p_{V_\psi}(\cdot),$$

$$k_p(x, y) = \nabla_x \nabla_y k(x, y) + \langle \nabla_x k(x, y), \nabla_y \log p(y) \rangle + \langle \nabla_y k(x, y), \nabla_x \log p(x) \rangle \\ + k(x, y) \langle \nabla_x \log p(x), \nabla_y \log p(y) \rangle$$

Kernel Stein discrepancy

Kernel Stein discrepancy between distributions (Riabiz et al. 2022)

$$KSD(p | p_{V_\psi}) = \sqrt{\frac{1}{n^2} \sum_{i,j=1}^n k_p(x_i, x_j)}, \quad x_i \sim p_{V_\psi}(\cdot),$$

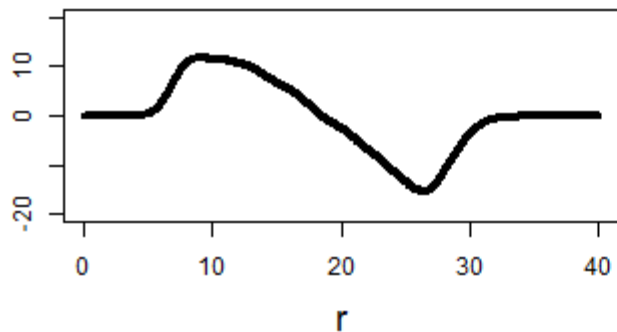
$$k_p(x, y) = \nabla_x \nabla_y k(x, y) + \langle \nabla_x k(x, y), \nabla_y \log p(y) \rangle + \langle \nabla_y k(x, y), \nabla_x \log p(x) \rangle \\ + k(x, y) \langle \nabla_x \log p(x), \nabla_y \log p(y) \rangle$$

Use kernelized Stein discrepancy for goodness-of-fit tests with $H_0 : p = p_{V_\psi}$

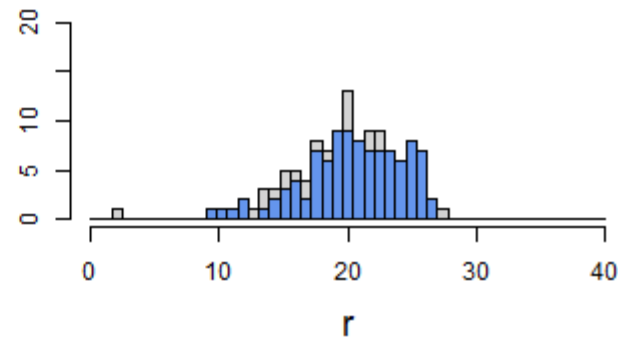
- Employing bootstrap procedure
- Stop when H_0 cannot be rejected anymore (significance level α)
- Conservative for correlated samples (Chwialkowski et al. 2016)

Stopping criteria

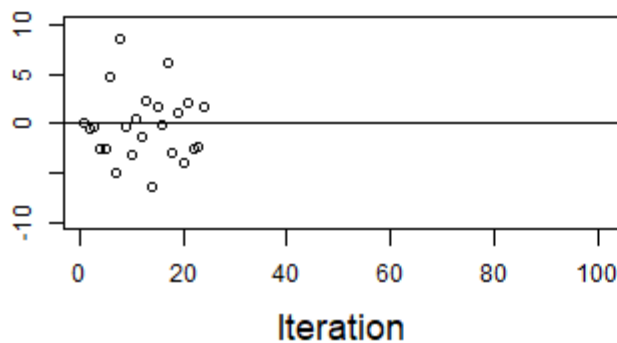
$V(r)$



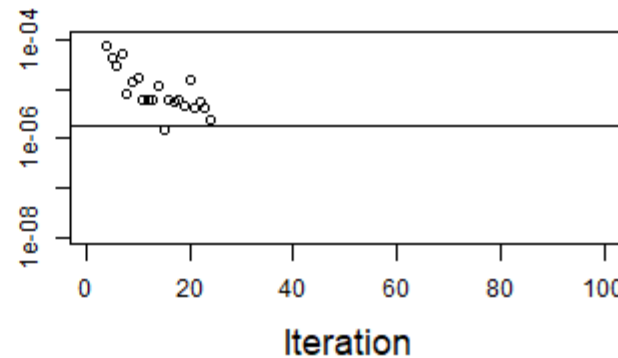
Samples $p(r)$ and $pV(r)$



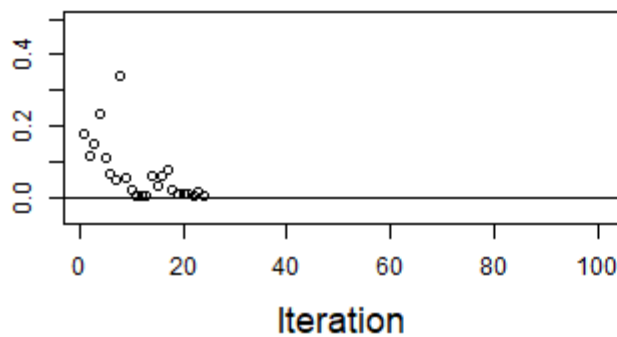
Loss



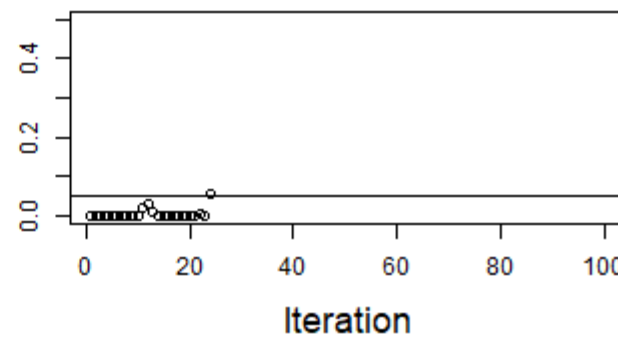
Risk probability



KSD

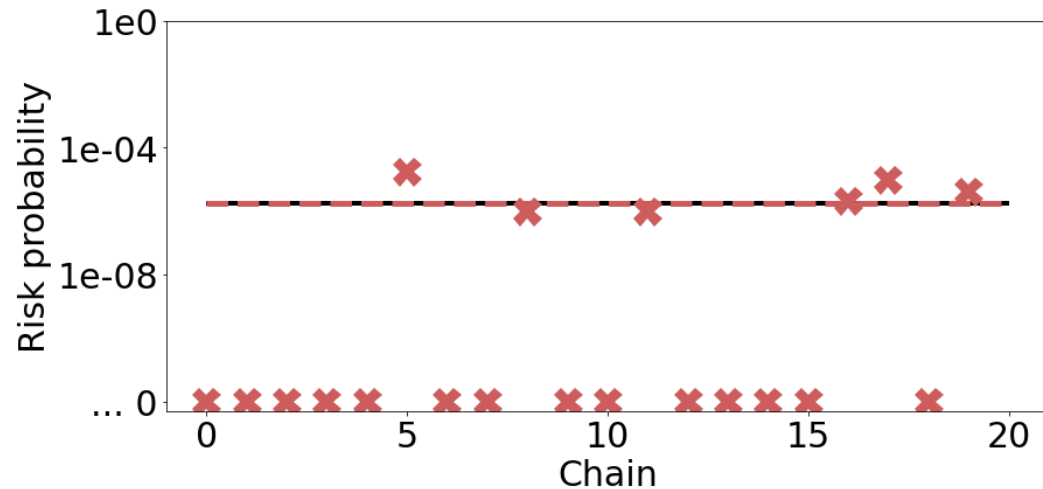


P-value



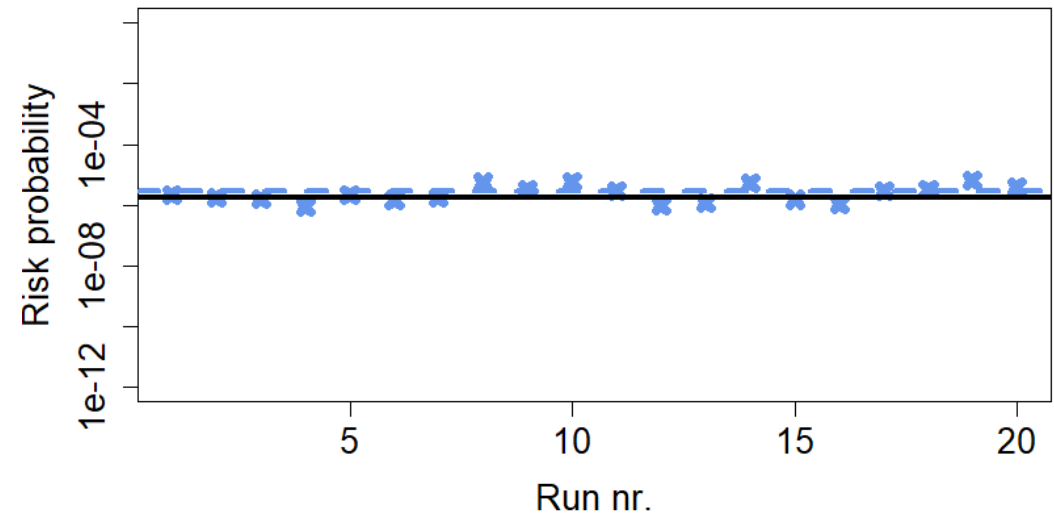
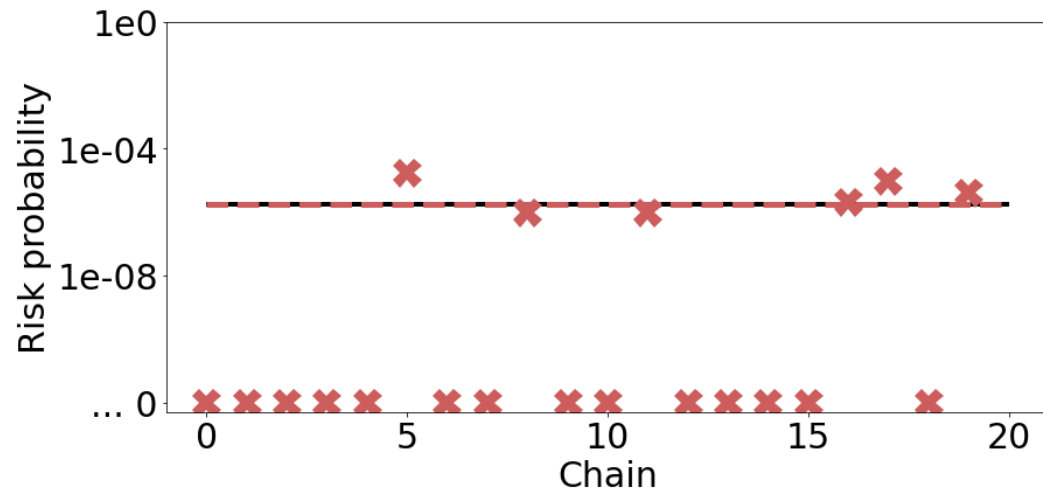
Energy based models vs. MCMC

- $\mathbb{P}(\mathcal{R}(\boldsymbol{\theta}) \geq 20.0 \mid \mathbf{y}) = 1.76 \times 10^{-6}$
- MCMC (left): 1'000'000 iterations \rightarrow SD = 4.2×10^{-6}
-



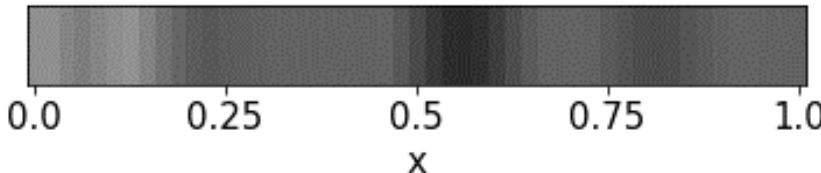
Energy based models vs. MCMC

- $\mathbb{P}(\mathcal{R}(\boldsymbol{\theta}) \geq 20.0 \mid \mathbf{y}) = 1.76 \times 10^{-6}$
- MCMC (left): 1'000'000 iterations \rightarrow SD = 4.2×10^{-6}
- EBM (right): 14'000-50'000 iterations \rightarrow SD = 1.5×10^{-6}



1D flow example (Straub et al. 2016)

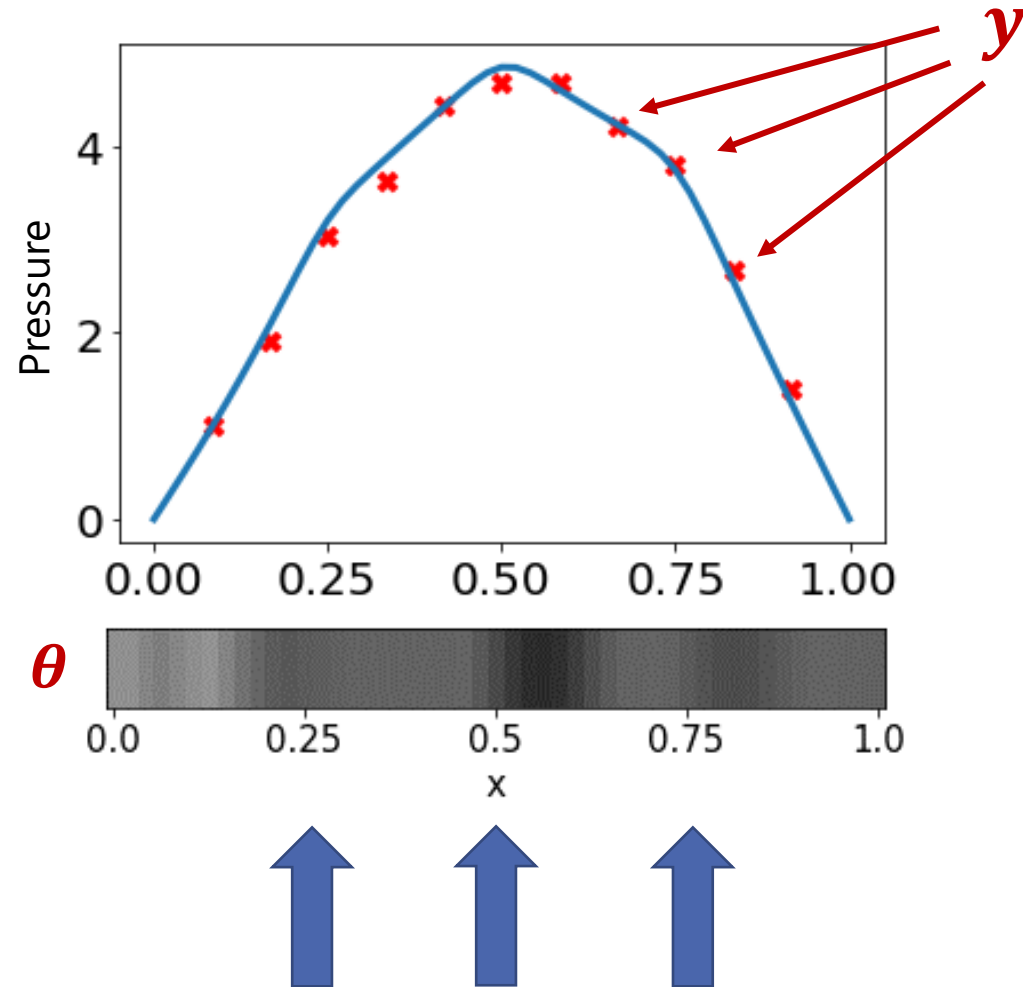
- Hydraulic diffusivity field $a(x)$, $x \in [0\text{m}, 1\text{m}]$
 - Log-diffusivity Gaussian, Karhunen-Loève expansion
- Diffusivity = speed at which pressure pulse propagates through aquifer

$$\ln a(x) = \mu_{\ln\theta} + \sigma_{\ln\theta} \sum_{i=1}^{10} \sqrt{w_i} v_i(x) Z_i$$


The figure shows a horizontal bar representing the spatial domain x from 0.0 to 1.0. The bar has a grayscale gradient, being darkest at $x=0.5$ and lighter at the ends. The x-axis is labeled with 0.0, 0.25, 0.5, 0.75, and 1.0. A red arrow points from the symbol θ to the sum in the equation above.

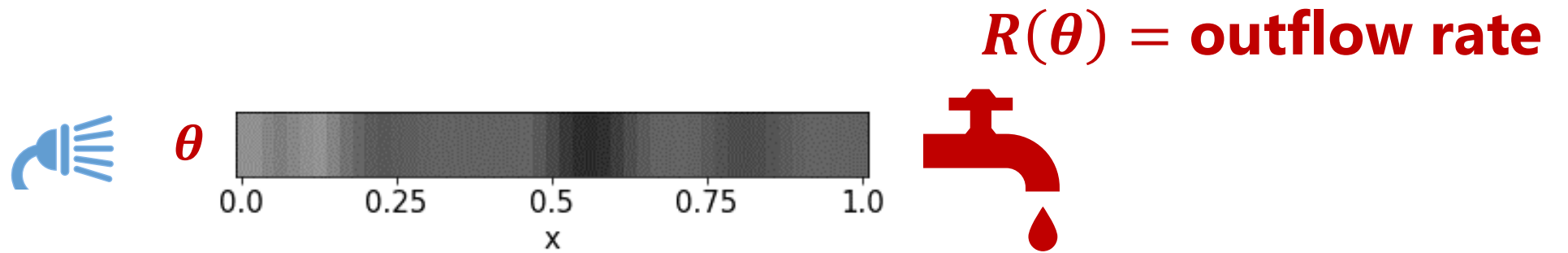
1D flow example

Data



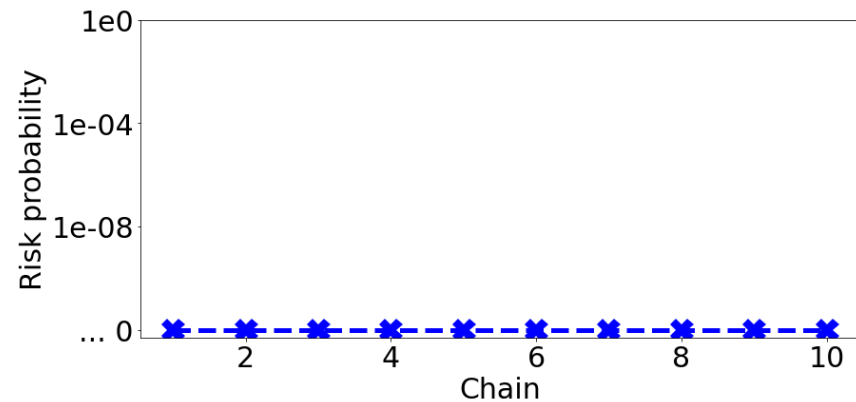
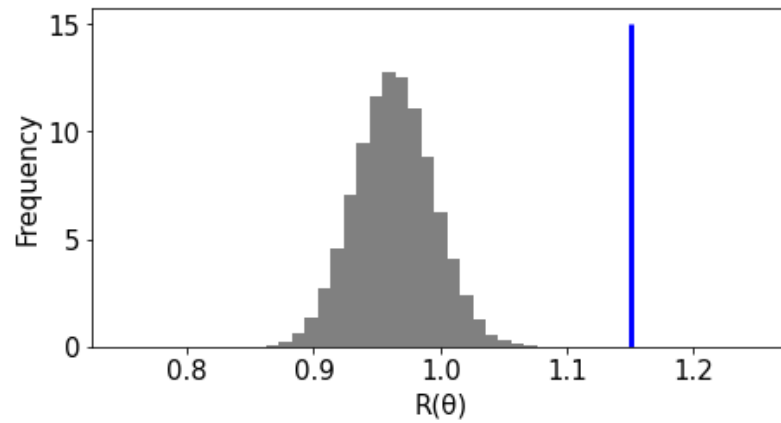
1D flow example

Risk



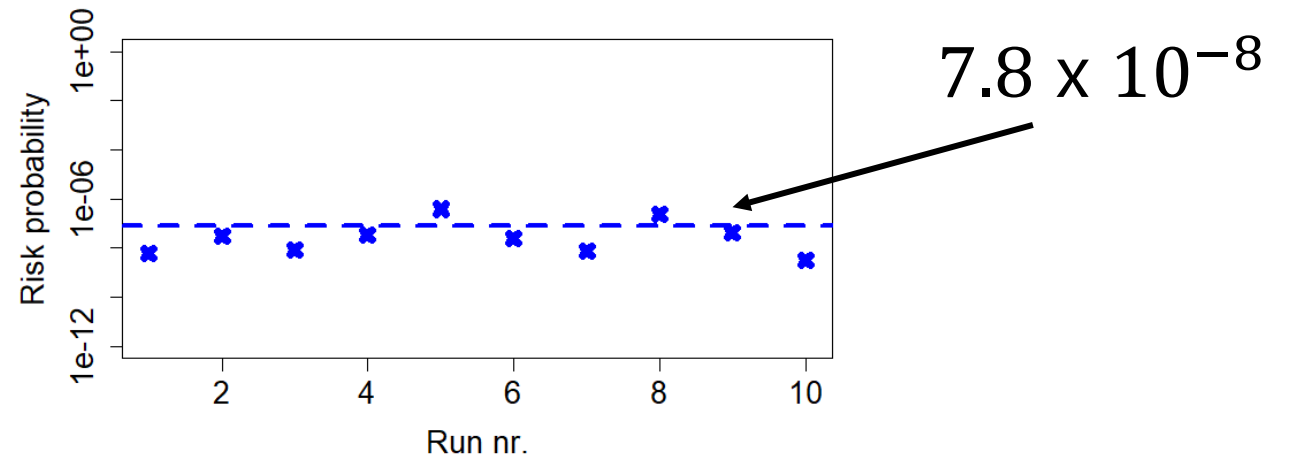
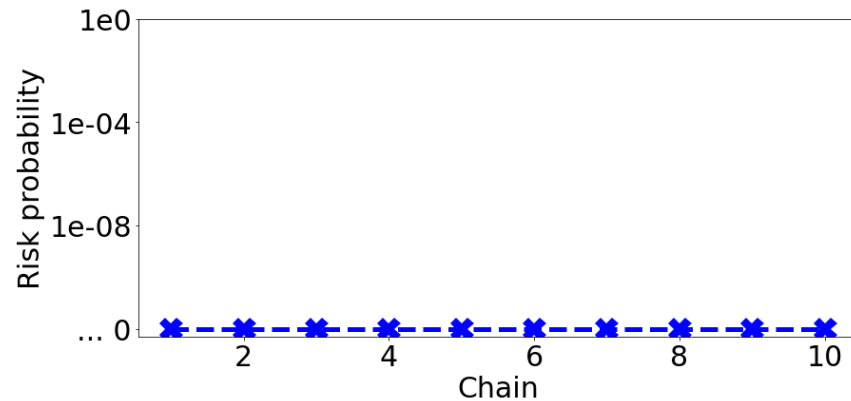
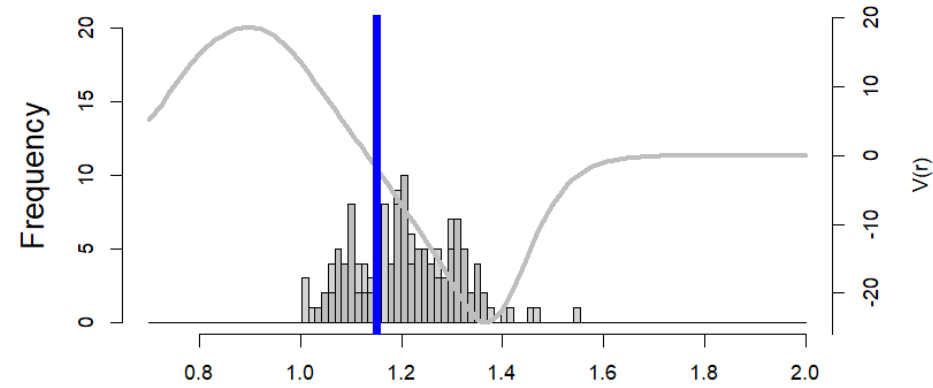
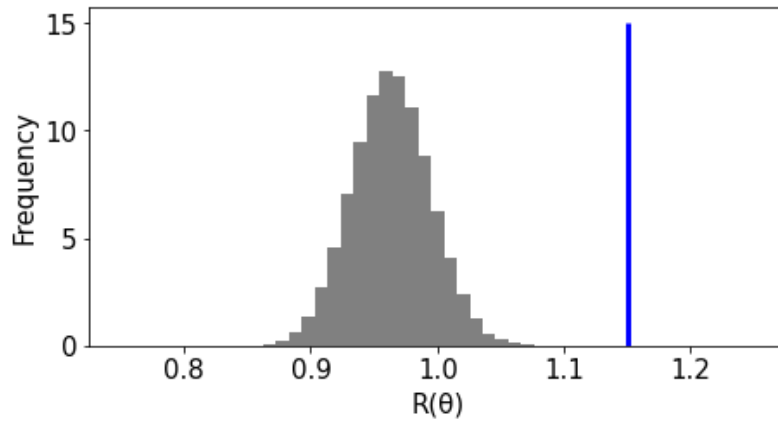
Energy based models vs. MCMC

$$\mathbb{P}(\mathcal{R}(\boldsymbol{\theta}) \geq 1.15 \mid \mathbf{y})$$



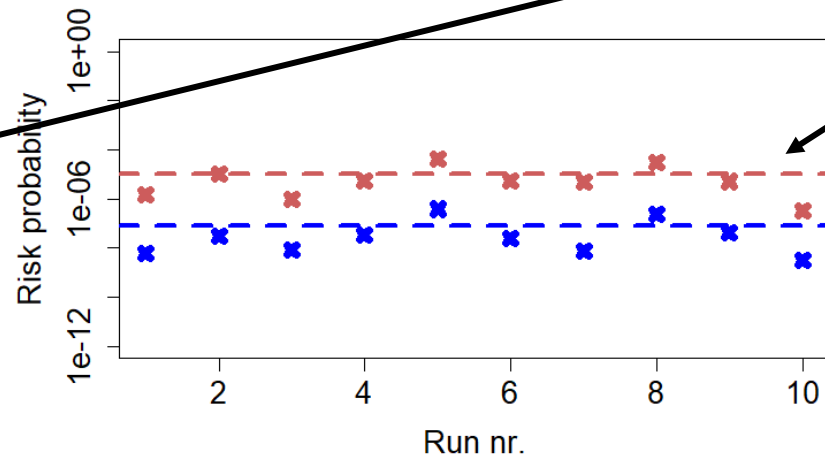
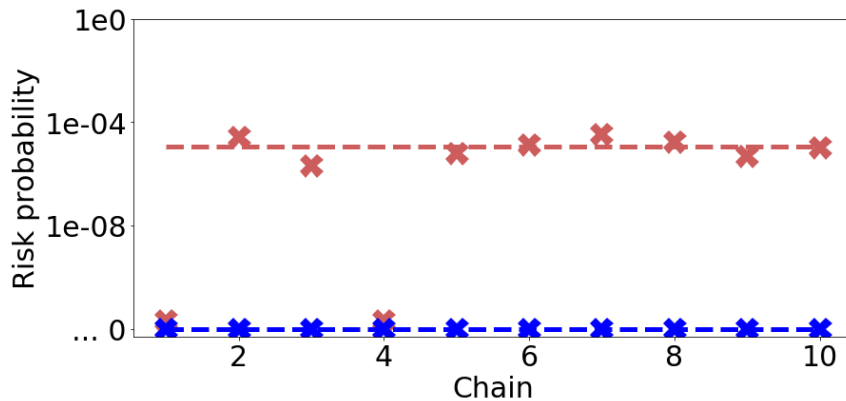
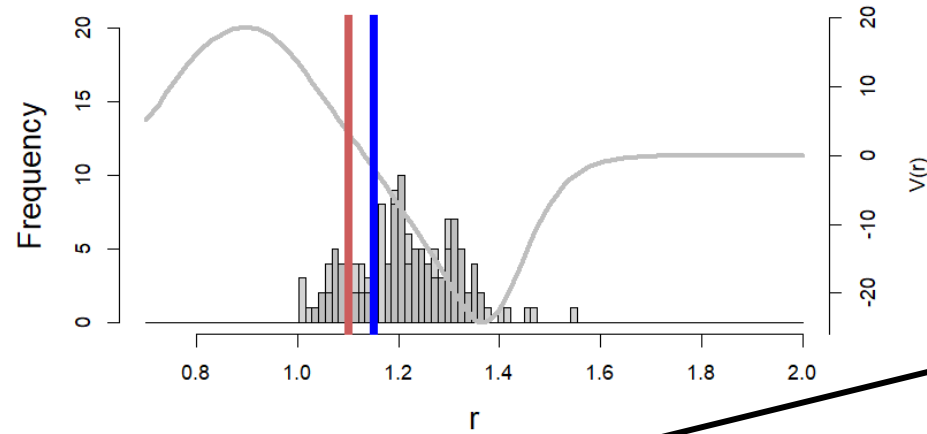
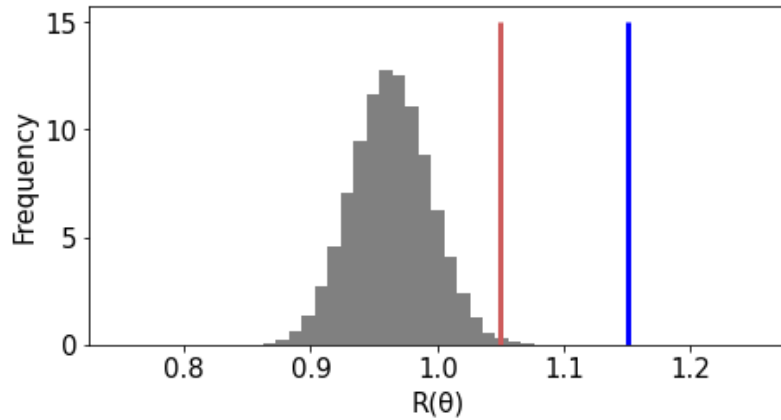
Energy based models vs. MCMC

$$\mathbb{P}(\mathcal{R}(\boldsymbol{\theta}) \geq 1.15 \mid \mathbf{y})$$



Energy based models vs. MCMC

$$+ \mathbb{P}(\mathcal{R}(\boldsymbol{\theta}) \geq 1.10 \mid \mathbf{y})$$



1.2×10^{-5}

1.1×10^{-5}

Conclusions and ongoing work

- **EBM approach reduces high-dimensional problem to optimization of one-dimensional function, compared to MCMC fraction of model evaluations needed**

Conclusions and ongoing work

- **EBM approach reduces high-dimensional problem to optimization of one-dimensional function, compared to MCMC fraction of model evaluations needed**
- **Configuration crucial**
 - Parameterization $V(r)$
 - Sampling $p_V(r)$
 - Choice of $p(r)$
 - Numerical integration
 - Learning rate
 - Stopping criteria

Conclusions and ongoing work

- **EBM approach reduces high-dimensional problem to optimization of one-dimensional function, compared to MCMC fraction of model evaluations needed**
- **Configuration crucial**
 - Parameterization $V(r)$
 - Sampling $p_V(r)$
 - Choice of $p(r)$
 - Numerical integration
 - Learning rate
 - Stopping criteria
- **Assessment and comparison**
 - MCMC to validate probabilities down to 10^{-6}
 - Reliability literature (Straub et al. 2016)
 - Sequential Monte Carlo



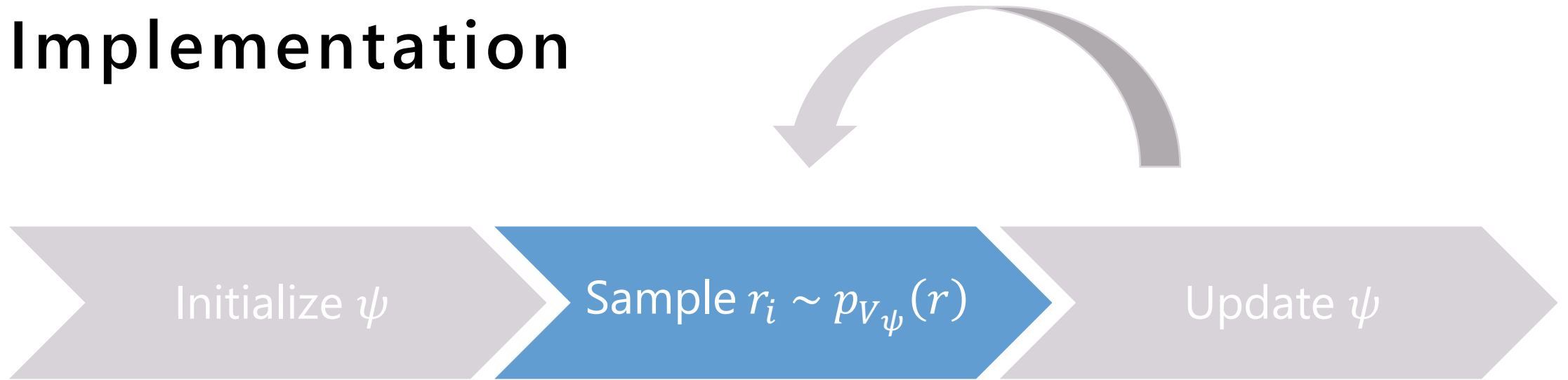
Thank you!

lea.friedli@unil.ch

References

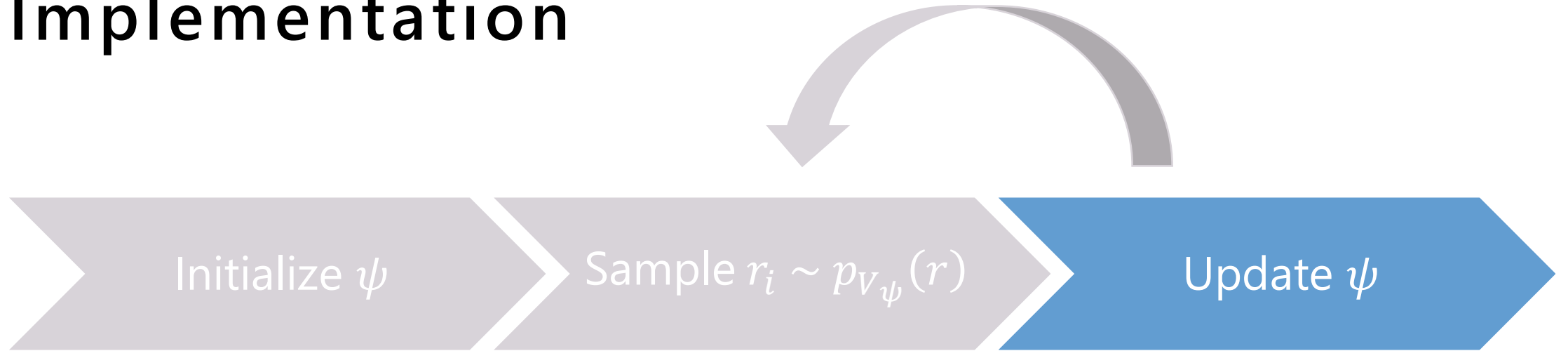
- Kacper Chwialkowski, Heiko Strathmann and Arthur Gretton. A kernel test of goodness of fit. In International conference on machine learning (pp. 2606-2615). PMLR, 2016.
- Michele Invernizzi, Omar Valsson, and Michele Parrinello. Coarse graining from variationally enhanced sampling applied to the ginzburg–landau model. *Proceedings of the National Academy of Sciences*, 114(13):3370–3374, 2017.
- Marina Riabiz, Wilson Ye Chen, Jon Cockayne, Pawel Swietach, Steven A. Niederer, Lester Mackey, and Chris. J. Oates. Optimal Thinning of MCMC Output. arXiv:2005.03952v5, 2022.
- Daniel Straub, Iason Papaioannou and Wolfgang Betz. Bayesian analysis of rare events. *Journal of Computational Physics*, 314, 538-556, 2016.
- Omar Valsson and Michele Parrinello. Variational approach to enhanced sampling and free energy calculations. *Physical Review Letters*, 113(9):090601, 2014.

Implementation



- $p_{V_\psi}(r) = \frac{\exp\left(-\left(F(r)+V_\psi(r)\right)\right)}{\int \exp\left(-\left(F(s)+V_\psi(s)\right)\right)ds}$
- Posterior PDF: $p(\theta|y) = \frac{\exp(-U(\theta))}{\int \exp(-U(\xi))d\xi}$ with $U(\theta) = -\log p(y|\theta) - \log p(\theta)$
- MCMC to draw proportional to $p_{V_\psi}(\theta) = \frac{\exp\left(-\left(U(\theta)+V_\psi(\mathcal{R}(\theta))\right)\right)}{\int \exp\left(-\left(U(\xi)+V_\psi(\mathcal{R}(\xi))\right)\right)d\xi}$
- Transform samples with $\theta \mapsto \mathcal{R}(\theta)$

Implementation



- Stochastic gradient descent $\psi_{n+1} = \psi_n - \gamma \frac{\partial J(\psi)}{\partial \psi}$
- Valsson and Parrinello (2014) employ loss related to KL divergence
- $\frac{\partial J(\psi)}{\partial \psi} \approx \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \psi} V_\psi(r_i) - \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \psi} V_\psi(s_i), s_i \sim p(\cdot), r_i \sim p_{V_\psi}(\cdot)$

Configuration analytical toy example

- **Parameterization $V(r)$**

1000 Gaussian radial basis functions (RBF, $\text{eps}=1$, from 0 to 40), learn weights

- **Learning rate**

0.5, geometric decrease with factor 1/1.025

- **MCMC:**

Metropolis-Hastings, Gaussian proposals, step width for AR 30%, 1100 steps, burn-in after 100 steps, thinning with factor 10

- **Choice of $p(r)$**

$\mathcal{N}(20,4)$

- **Stopping criteria**

$\alpha = 0.05$, 1000 bootstrap samples

Configuration 1D flow example

- **Parameterization $V(r)$**

1000 Gaussian radial basis functions (RBF, $\text{eps}=20$ from 0 to 2), learn weights

- **Learning rate**

0.25, geometric decrease with factor 1/1.025

- **MCMC:**

Metropolis-Hastings, Gaussian proposals, step width for AR 30%, 1200 steps, burn-in after 200 steps, thinning with factor 10

- **Choice of $p(r)$**

$\mathcal{N}(1.2, 0.125)$

- **Stopping criteria**

$\alpha = 0.01$, 1000 bootstrap samples