

# Fast and Correct Gradient-Based Optimisation for Probabilistic Programming via Smoothing

Basim Khajwal   Luke Ong   Dominik Wagner



32th European Symposium on Programming  
26 April 2023

# Probabilistic Programming

# Probabilistic Programming

= programming paradigm to pose *Bayesian Inference* problems

# Probabilistic Programming

= programming paradigm to pose *Bayesian Inference* problems

- ▶ separate modelling from inference

# Variational Inference

## Variational Inference:

*frame posterior inference as (deterministic) optimisation problem*

# Variational Inference:

*frame posterior inference as (deterministic) optimisation problem*

**Posit:** *variational family* of “simpler” *guide* distributions

# Variational Inference:

*frame posterior inference as (deterministic) optimisation problem*

**Posit:** *variational family* of “simpler” *guide* distributions

**Aim:** find guide that is “*closest*” to (true) posterior



# Variational Inference:

frame posterior inference as (deterministic) *optimisation* problem

**Posit:** *variational family* of “simpler” *guide* distributions

**Aim:** find guide that is “*closest*” to (true) posterior  
 KL divergence

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}_{\theta}} [f(\theta, \mathbf{s})]$$

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}_{\theta}} [f(\theta, \mathbf{s})]$$

use *Stochastic Gradient Descent*

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}_{\theta}} [f(\theta, \mathbf{s})]$$

use *Stochastic Gradient Descent*

**Key ingredient:** estimation of gradient of expectation

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}_{\theta}} [f(\theta, \mathbf{s})]$$

use *Stochastic Gradient Descent*

**Key ingredient:** estimation of gradient of expectation

- **Score Estimator**

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}_{\theta}} [f(\theta, \mathbf{s})]$$

use *Stochastic Gradient Descent*

**Key ingredient:** estimation of gradient of expectation

- **Score Estimator:**  
widely applicable but *high variance*

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}_{\theta}} [f(\theta, \mathbf{s})]$$

use *Stochastic Gradient Descent*

**Key ingredient:** estimation of gradient of expectation

- **Score Estimator:**  
widely applicable but *high variance*
- **Reparameterisation Estimator**

# Reparameterisation Gradient Estimator

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}_{\theta}} [f(\theta, \mathbf{s})]$$



# Reparameterisation Gradient Estimator

$$\operatorname{argmin}_{\theta} \mathbb{E}_{s \sim D_{\theta}} [f(\theta, s)]$$

eliminate dependence on  $\theta$



# Reparameterisation Gradient Estimator

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [f(\theta, \mathbf{s})]$$

eliminate dependence on  $\theta$

estimate:  $\nabla_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [f(\theta, \mathbf{s})] \approx \nabla_{\theta} f(\theta, \hat{\mathbf{s}})$ , where  $\hat{\mathbf{s}} \sim \mathcal{D}$

# Reparameterisation Gradient Estimator

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [f(\theta, \mathbf{s})]$$

eliminate dependence on  $\theta$

estimate:  $\nabla_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [f(\theta, \mathbf{s})] \approx \nabla_{\theta} f(\theta, \hat{\mathbf{s}})$ , where  $\hat{\mathbf{s}} \sim \mathcal{D}$

**(Unbiasedness)**  $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [\nabla_{\theta} f(\theta, \mathbf{s})] \stackrel{?}{=} \nabla_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [f(\theta, \mathbf{s})]$

# Reparameterisation Gradient Estimator

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [f(\theta, \mathbf{s})]$$

eliminate dependence on  $\theta$

expressed in PL with conditionals,  
may *not* be *differentiable/continuous*

estimate:  $\nabla_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [f(\theta, \mathbf{s})] \approx \nabla_{\theta} f(\theta, \hat{\mathbf{s}})$ , where  $\hat{\mathbf{s}} \sim \mathcal{D}$

(Unbiasedness)  $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [\nabla_{\theta} f(\theta, \mathbf{s})] \stackrel{?}{=} \nabla_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [f(\theta, \mathbf{s})]$

# Reparameterisation Gradient Estimator

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [f(\theta, \mathbf{s})]$$

eliminate dependence on  $\theta$

expressed in PL with conditionals,  
may *not* be *differentiable/continuous*

estimate:  $\nabla_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [f(\theta, \mathbf{s})] \approx \nabla_{\theta} f(\theta, \hat{\mathbf{s}})$ , where  $\hat{\mathbf{s}} \sim \mathcal{D}$

(Unbiasedness)

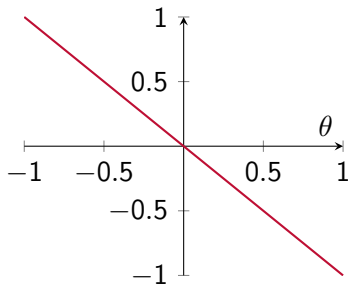
$$\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [\nabla_{\theta} f(\theta, \mathbf{s})] \stackrel{?}{=} \nabla_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [f(\theta, \mathbf{s})]$$

*may be compromised!* [Lee et al., NeurIPS 2018]

$$f(\theta, s) = -0.5 \cdot \theta^2 + \begin{cases} 0 & \text{if } s + \theta < 0 \\ 1 & \text{otherwise} \end{cases}$$

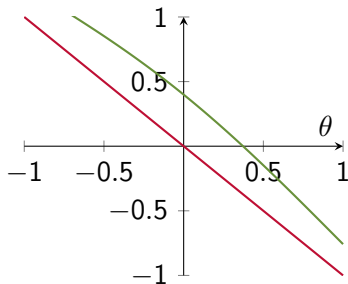
$$f(\theta, s) = -0.5 \cdot \theta^2 + \begin{cases} 0 & \text{if } s + \theta < 0 \\ 1 & \text{otherwise} \end{cases}$$

$$\mathbb{E}_{\mathbf{s} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})} [\nabla_{\theta} \mathbf{f}(\theta, \mathbf{s})] = -\theta$$



$$f(\theta, s) = -0.5 \cdot \theta^2 + \begin{cases} 0 & \text{if } s + \theta < 0 \\ 1 & \text{otherwise} \end{cases}$$

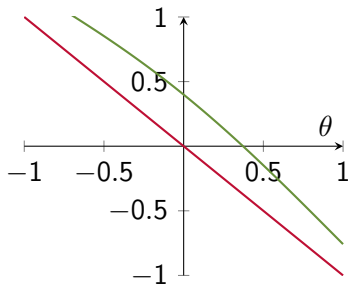
$$\mathbb{E}_{s \sim \mathcal{N}(0,1)} [\nabla_{\theta} f(\theta, s)] = -\theta \neq -\theta + \mathcal{N}(-\theta | 0, 1) = \nabla_{\theta} \mathbb{E}_{s \sim \mathcal{N}(0,1)} [f(\theta, s)]$$





$$f(\theta, s) = -0.5 \cdot \theta^2 + \begin{cases} 0 & \text{if } s + \theta < 0 \\ 1 & \text{otherwise} \end{cases}$$

$$\mathbb{E}_{s \sim \mathcal{N}(0,1)} [\nabla_{\theta} f(\theta, s)] = -\theta \neq -\theta + \mathcal{N}(-\theta | 0, 1) = \nabla_{\theta} \mathbb{E}_{s \sim \mathcal{N}(0,1)} [f(\theta, s)]$$



**Stochastic Gradient Descent is *incorrect!***

# Contributions

*Fast yet correct* Stochastic Gradient Descent  
with Reparameterisation Gradient *via Smoothing*

# Contributions

*Fast yet correct* Stochastic Gradient Descent  
with Reparameterisation Gradient *via Smoothing*

- ▶ Smoothed Denotational (Value) Semantics

# Contributions

*Fast yet correct* Stochastic Gradient Descent  
with Reparameterisation Gradient *via Smoothing*

- ▶ Smoothed Denotational (Value) Semantics
- ▶ Correctness of Stochastic Gradient Descent via Type System

# Contributions

*Fast yet correct* Stochastic Gradient Descent  
with Reparameterisation Gradient *via Smoothing*

- ▶ Smoothed Denotational (Value) Semantics
- ▶ Correctness of Stochastic Gradient Descent via Type System
- ▶ Convergence of Smooth Approximations

# Contributions

*Fast yet correct* Stochastic Gradient Descent  
with Reparameterisation Gradient *via Smoothing*

- ▶ Smoothed Denotational (Value) Semantics
- ▶ Correctness of Stochastic Gradient Descent via Type System
- ▶ Convergence of Smooth Approximations
- ▶ Empirical Evaluation

Part I:  
Problem Setup

simply typed  $\lambda$ -calculus with  $\mathbb{R}$ , primitive operations, parameters  $\theta_i$

$$M ::= x \mid \lambda x. M \mid M M \mid f(M, \dots, M) \mid \theta_i$$



simply typed  $\lambda$ -calculus with  $\mathbb{R}$ , primitive operations, parameters  $\theta_i$   
+ sample

$$M ::= x \mid \lambda x. M \mid M M \mid f(M, \dots, M) \mid \theta_i \\ \mid \mathbf{sample}_{\mathcal{D}}$$

simply typed  $\lambda$ -calculus with  $\mathbb{R}$ , primitive operations, parameters  $\theta_i$

+ sample

+ *branching*

$$M ::= x \mid \lambda x. M \mid M M \mid f(M, \dots, M) \mid \theta_i$$

| **sample** <sub>$\mathcal{D}$</sub>

| **if**  $M < 0$  **then**  $M$  **else**  $M$

## Denotational Value Semantics:

Denotational Value Semantics: *deterministic* function from samples to value

Denotational Value Semantics: *deterministic* function from samples to value

$$f(\theta, s) = -0.5 \cdot \theta^2 + \begin{cases} 0 & \text{if } s + \theta < 0 \\ 1 & \text{otherwise} \end{cases}$$

Denotational Value Semantics: *deterministic* function from samples to value

$$f(\theta, s) = -0.5 \cdot \theta^2 + \begin{cases} 0 & \text{if } s + \theta < 0 \\ 1 & \text{otherwise} \end{cases}$$

$$M \equiv (\lambda s. -0.5 \cdot \theta^2 + (\text{if } s + \theta < 0 \text{ then } 0 \text{ else } 1)) \text{sample}_{\mathcal{N}}$$

Denotational Value Semantics: *deterministic* function from samples to value

$$\llbracket M \rrbracket (\theta, s) = f(\theta, s) = -0.5 \cdot \theta^2 + \begin{cases} 0 & \text{if } s + \theta < 0 \\ 1 & \text{otherwise} \end{cases}$$

$$M \equiv (\lambda s. -0.5 \cdot \theta^2 + (\text{if } s + \theta < 0 \text{ then } 0 \text{ else } 1)) \text{sample}_{\mathcal{N}}$$

Denotational Value Semantics: *deterministic* function from samples to value

$$\llbracket M \rrbracket (\theta, s) = f(\theta, s) = -0.5 \cdot \theta^2 + \begin{cases} 0 & \text{if } s + \theta < 0 \\ 1 & \text{otherwise} \end{cases}$$

$$M \equiv (\lambda s. -0.5 \cdot \theta^2 + (\text{if } s + \theta < 0 \text{ then } 0 \text{ else } 1)) \text{sample}_{\mathcal{N}}$$

*Track samples (and distributions) in type system*



Denotational Value Semantics: *deterministic* function from samples to value

$$\llbracket M \rrbracket (\theta, s) = f(\theta, s) = -0.5 \cdot \theta^2 + \begin{cases} 0 & \text{if } s + \theta < 0 \\ 1 & \text{otherwise} \end{cases}$$

$$M \equiv (\lambda s. -0.5 \cdot \theta^2 + (\text{if } s + \theta < 0 \text{ then } 0 \text{ else } 1)) \text{sample}_{\mathcal{N}}$$

*Track samples (and distributions) in type system*

$$\theta : R \mid [\mathcal{N}] \vdash M : R$$

# Problem Statement

**Given:** term-in-context,  $\theta_1 : R, \dots, \theta_m : R \mid [\mathcal{D}_1, \dots, \mathcal{D}_n] \vdash M : R$

# Problem Statement

**Given:** term-in-context,  $\theta_1 : R, \dots, \theta_m : R \mid [\mathcal{D}_1, \dots, \mathcal{D}_n] \vdash M : R$

**Find:**  $\operatorname{argmin}_{\theta} \mathbb{E}_{s_1 \sim \mathcal{D}_1, \dots, s_n \sim \mathcal{D}_n} [\llbracket M \rrbracket (\theta, \mathbf{s})]$

$$\mathbb{E}_{s \sim \mathcal{N}} [\exp(s^2)] = \infty$$

$$\mathbb{E}_{s \sim \mathcal{N}} [\exp(s^2)] = \infty$$

$(\lambda x. \exp(x \cdot x)) \text{ sample}_{\mathcal{N}}$

$$\mathbb{E}_{s \sim \mathcal{N}} [\exp(s^2)] = \infty$$

$(\lambda x. \exp(x \cdot x)) \text{ sample}_{\mathcal{N}}$

**Simplified assumption:**

1. distributions have finite moments

$$\mathbb{E}_{s \sim \mathcal{N}} [\exp(s^2)] = \infty$$

$(\lambda x. \exp(x \cdot x)) \text{ sample}_{\mathcal{N}}$

**Simplified assumption:**

1. distributions have finite moments

$$\mathbb{E}_{s \sim \mathcal{D}} [|s^p|] < \infty$$

$$\mathbb{E}_{s \sim \mathcal{N}} [\exp(s^2)] = \infty$$

$(\lambda x. \exp(x \cdot x)) \text{ sample}_{\mathcal{N}}$

### Simplified assumption:

1. distributions have finite moments
2. primitives are bounded by polynomials

$$\mathbb{E}_{s \sim \mathcal{D}} [|s^p|] < \infty$$



$$\mathbb{E}_{s \sim \mathcal{N}} [\exp(s^2)] = \infty$$

$(\lambda x. \exp(x \cdot x)) \text{ sample}_{\mathcal{N}}$

### Simplified assumption:

1. distributions have finite moments
2. primitives are bounded by polynomials

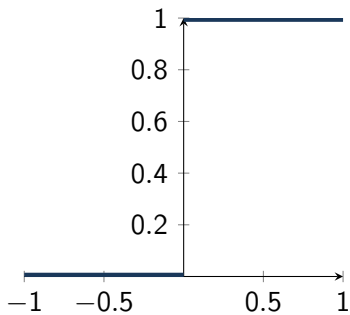
$$\mathbb{E}_{s \sim \mathcal{D}} [|s^p|] < \infty$$

**In paper:** relax assumption, control use of  $\log$ ,  $\exp$ ,  $^{-1}$  via type system

Part II:  
Smoothed Value Semantics

$$\llbracket \text{if } z < 0 \text{ then } 0 \text{ else } M \rrbracket (z) = \begin{cases} 0 & \text{if } z < 0 \\ \llbracket M \rrbracket (z) & \text{otherwise} \end{cases}$$

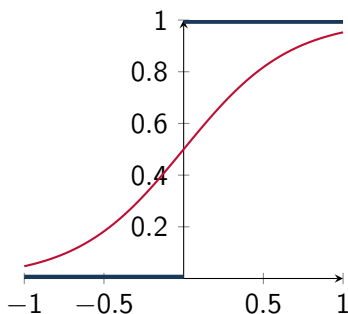
$$\llbracket \text{if } z < 0 \text{ then } 0 \text{ else } M \rrbracket (z) = [z \geq 0] \cdot \llbracket M \rrbracket (z)$$



$$\llbracket \text{if } z < 0 \text{ then } 0 \text{ else } M \rrbracket (z) = [z \geq 0] \cdot \llbracket M \rrbracket (z)$$

$$\llbracket \text{if } z < 0 \text{ then } 0 \text{ else } M \rrbracket_{\eta} (z) = \sigma_{\eta}(z) \cdot \llbracket M \rrbracket_{\eta} (z)$$

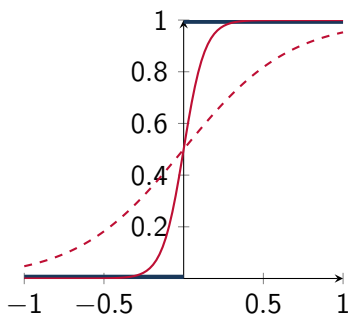
sigmoid function



$$\llbracket \text{if } z < 0 \text{ then } 0 \text{ else } M \rrbracket (z) = [z \geq 0] \cdot \llbracket M \rrbracket (z)$$

$$\llbracket \text{if } z < 0 \text{ then } 0 \text{ else } M \rrbracket_{\eta} (z) = \sigma_{\eta}(z) \cdot \llbracket M \rrbracket_{\eta} (z)$$

sigmoid function (parameterised by *accuracy* coefficient  $\eta > 0$ )



# Categorical Model?

# Categorical Model?

Smoothness: CCC of *Frölicher* spaces



# Categorical Model?

Smoothness: CCC of *Frölicher* spaces

$$\llbracket \text{if } L < 0 \text{ then } M \text{ else } N \rrbracket_{\eta} := (\sigma_{\eta} \circ (- \llbracket L \rrbracket_{\eta})) \cdot \llbracket M \rrbracket_{\eta} + (\sigma_{\eta} \circ \llbracket L \rrbracket_{\eta}) \cdot \llbracket N \rrbracket_{\eta}$$

# Categorical Model?

Smoothness: CCC of *Frölicher* spaces

$$\llbracket \text{if } L < 0 \text{ then } M \text{ else } N \rrbracket_{\eta} := (\sigma_{\eta} \circ (- \llbracket L \rrbracket_{\eta})) \cdot \llbracket M \rrbracket_{\eta} + (\sigma_{\eta} \circ \llbracket L \rrbracket_{\eta}) \cdot \llbracket N \rrbracket_{\eta}$$

*morphism?*

# Categorical Model?

Smoothness: CCC of *Frölicher* spaces

$$\llbracket \text{if } L < 0 \text{ then } M \text{ else } N \rrbracket_{\eta} := (\sigma_{\eta} \circ (- \llbracket L \rrbracket_{\eta})) \cdot \llbracket M \rrbracket_{\eta} + (\sigma_{\eta} \circ \llbracket L \rrbracket_{\eta}) \cdot \llbracket N \rrbracket_{\eta}$$

*morphism?*

*Adapt construction of Frölicher spaces*

# Categorical Model?

Smoothness: CCC of *Frölicher* spaces

$$\llbracket \text{if } L < 0 \text{ then } M \text{ else } N \rrbracket_{\eta} := (\sigma_{\eta} \circ (- \llbracket L \rrbracket_{\eta})) \cdot \llbracket M \rrbracket_{\eta} + (\sigma_{\eta} \circ \llbracket L \rrbracket_{\eta}) \cdot \llbracket N \rrbracket_{\eta}$$

*morphism?*

*Adapt construction of Frölicher spaces*

- + vector space structure for underlying set
- + condition for “curves”

# Categorical Model?

Smoothness: CCC of *Frölicher* spaces

$$\llbracket \text{if } L < 0 \text{ then } M \text{ else } N \rrbracket_{\eta} := (\sigma_{\eta} \circ (- \llbracket L \rrbracket_{\eta})) \cdot \llbracket M \rrbracket_{\eta} + (\sigma_{\eta} \circ \llbracket L \rrbracket_{\eta}) \cdot \llbracket N \rrbracket_{\eta}$$

*morphism?*

*Adapt construction of Frölicher spaces*

- + vector space structure for underlying set
- + condition for “curves”

CCC **VectFr** of *Vector Frölicher Spaces*

# Categorical Model?

Smoothness: CCC of *Frölicher* spaces

$$\llbracket \text{if } L < 0 \text{ then } M \text{ else } N \rrbracket_{\eta} := (\sigma_{\eta} \circ (- \llbracket L \rrbracket_{\eta})) \cdot \llbracket M \rrbracket_{\eta} + (\sigma_{\eta} \circ \llbracket L \rrbracket_{\eta}) \cdot \llbracket N \rrbracket_{\eta}$$

*morphism?*

*Adapt construction of Frölicher spaces*

- + vector space structure for underlying set
- + condition for “curves”

CCC **VectFr** of *Vector Frölicher Spaces*

If  $\phi_1, \phi_2 \in \mathbf{VectFr}(X, Y)$  and  $\alpha \in \mathbf{Vect}(X, \mathbb{R})$  then  $\alpha \cdot \phi_1 + \phi_2 \in \mathbf{VectFr}(X, Y)$ .

Part III:  
Applying Stochastic Gradient  
Descent

$$\boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k - \gamma_k \cdot \nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}_k, \mathbf{s}_k)$$

$$\mathbf{s}_k \sim \mathcal{D}$$



$$\boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k - \gamma_k \cdot \underbrace{\nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}_k, \mathbf{s}_k)}_{\text{gradient estimation}}$$

$$\mathbf{s}_k \sim \mathcal{D}$$

$$\boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k - \underbrace{\gamma_k \cdot \nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}_k, \mathbf{s}_k)}_{\text{gradient estimation}}$$

step size

$$\mathbf{s}_k \sim \mathcal{D}$$

$$\boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k - \underbrace{\gamma_k \cdot \nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}_k, \mathbf{s}_k)}_{\text{gradient estimation}} \quad \mathbf{s}_k \sim \mathcal{D}$$

step size

**(Unbiasedness)**  $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})] = \nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})]$

$$\boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k - \underbrace{\gamma_k \cdot \nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}_k, \mathbf{s}_k)}_{\text{gradient estimation}} \quad \mathbf{s}_k \sim \mathcal{D}$$

step size

**(Unbiasedness)**  $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})] = \nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})]$

*partial derivatives of  $\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})$  are bounded by polynomial*

$$\boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k - \underbrace{\gamma_k \cdot \nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}_k, \mathbf{s}_k)}_{\text{gradient estimation}} \quad \mathbf{s}_k \sim \mathcal{D}$$

step size

**(Unbiasedness)**  $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})] = \nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})]$

*partial derivatives of  $\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})$  are bounded by polynomial*

## Correctness of SGD for Smoothing

If  $M$  is typable,  $\Theta$  is compact and the step size scheme is “suitable”

$$\boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k - \underbrace{\gamma_k \cdot \nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}_k, \mathbf{s}_k)}_{\text{gradient estimation}} \quad \mathbf{s}_k \sim \mathcal{D}$$

step size

**(Unbiasedness)**  $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})] = \nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})]$

*partial derivatives of  $\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})$  are bounded by polynomial*

## Correctness of SGD for Smoothing

If  $M$  is typable,  $\Theta$  is compact and the step size scheme is “suitable” then

$$\inf_{i \in \mathbb{N}} \mathbb{E}[\nabla g(\boldsymbol{\theta}_i)] = 0$$

where  $g(\boldsymbol{\theta}) := \mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})]$ .

$$\boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k - \underbrace{\gamma_k \cdot \nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}_k, \mathbf{s}_k)}_{\text{gradient estimation}} \quad \mathbf{s}_k \sim \mathcal{D}$$

step size

**(Unbiasedness)**  $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\nabla_{\boldsymbol{\theta}} \llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})] = \nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})]$

*partial derivatives of  $\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})$  are bounded by polynomial*

## Correctness of SGD for Smoothing

If  $M$  is typable,  $\Theta$  is compact and the step size scheme is “suitable” then

$$\inf_{i \in \mathbb{N}} \mathbb{E}[\nabla g(\boldsymbol{\theta}_i)] = 0$$

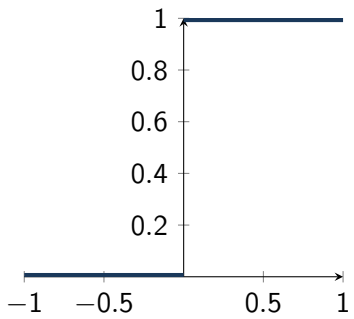
where  $g(\boldsymbol{\theta}) := \mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})]$ .

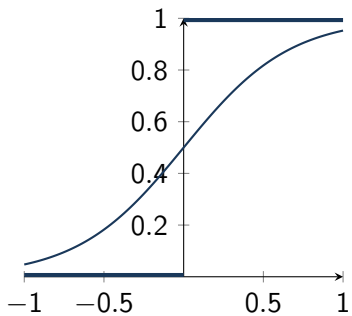
*exploit Lipschitz smoothness and bounded variance*

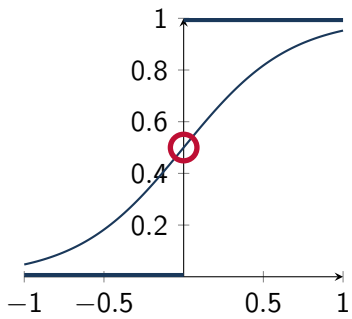
*How does solving the **smoothed** problem help solve the **original** problem?*

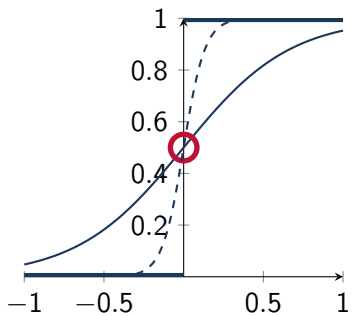


Part IV:  
Convergence of Smoothings



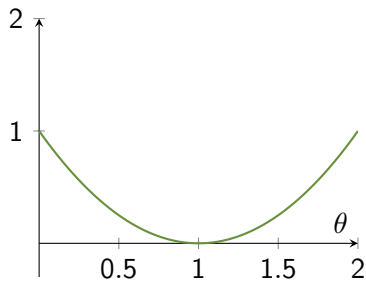






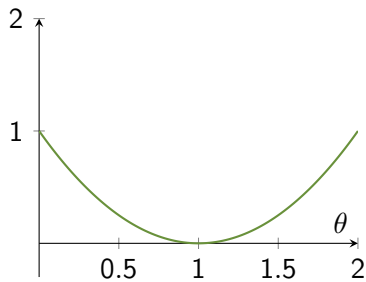
$$M \equiv \text{if } 0 < 0 \text{ then } \theta^2 + 1 \text{ else } (\theta - 1)^2$$

$$M \equiv \text{if } 0 < \theta \text{ then } \theta^2 + 1 \text{ else } (\theta - 1)^2$$



$$M \equiv \text{if } 0 < \theta \text{ then } \theta^2 + 1 \text{ else } (\theta - 1)^2$$

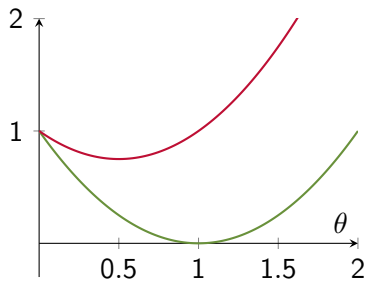
$$\llbracket M \rrbracket_{\eta}(\theta) = \frac{1}{2}(\theta^2 + 1) + \frac{1}{2}(\theta - 1)^2$$





$$M \equiv \text{if } 0 < \theta \text{ then } \theta^2 + 1 \text{ else } (\theta - 1)^2$$

$$\llbracket M \rrbracket_{\eta}(\theta) = \frac{1}{2}(\theta^2 + 1) + \frac{1}{2}(\theta - 1)^2$$



*Ensure that guards are **not 0 almost everywhere***

if  $x - x < 0$  then  $M$  else  $N$  **X**

**if  $x - x < 0$  then  $M$  else  $N$   $\times$**

**$(\lambda y, z. \text{if } y - z < 0 \text{ then } M \text{ else } N)_{xx}$   $\times$**

if  $x - x < 0$  then  $M$  else  $N$  **X**

$(\lambda y, z. \text{if } y - z < 0 \text{ then } M \text{ else } N) x x$  **X**

if  $\theta < 0$  then  $M$  else  $N$  **X**

if  $x - x < 0$  then  $M$  else  $N$  ✗

$(\lambda y, z. \text{if } y - z < 0 \text{ then } M \text{ else } N)$   $x$   $x$  ✗

if  $\theta < 0$  then  $M$  else  $N$  ✗

if  $(\text{sample}_{\mathcal{N}} + \theta) < 0$  then  $M$  else  $N$  ✓

if  $x - x < 0$  then  $M$  else  $N$  ✗

$(\lambda y, z. \text{if } y - z < 0 \text{ then } M \text{ else } N) x x$  ✗

if  $\theta < 0$  then  $M$  else  $N$  ✗

if  $(\underbrace{\text{sample}_{\mathcal{N}} + \theta}) < 0$  then  $M$  else  $N$  ✓  
transform  $\text{sample}_{\mathcal{N}}$  by  $(\lambda x. x + \theta)$

if  $x - x < 0$  then  $M$  else  $N$  ✗

$(\lambda y, z. \text{if } y - z < 0 \text{ then } M \text{ else } N) x x$  ✗

if  $\theta < 0$  then  $M$  else  $N$  ✗

if  $(\underbrace{\text{sample}_{\mathcal{N}} + \theta}) < 0$  then  $M$  else  $N$  ✓  
transform  $\text{sample}_{\mathcal{N}}$  by  $(\lambda x. x + \theta)$

$(\lambda y, z. \text{if } y - z < 0 \text{ then } M \text{ else } N) \text{sample}_{\mathcal{N}} (\text{transform } \text{sample}_{\mathcal{N}} \text{ by } T)$  ✓



$$\tau ::= R(g, \Delta)$$

$$\tau ::= R(g, \Delta)$$

guard-safe?



$$\tau ::= R(g, \Delta)$$

guard-safe?

dependency on (transformed) samples

$$\tau ::= R(g, \Delta) \mid \tau \rightarrow \tau$$

guard-safe?

dependency on (transformed) samples

$$\tau ::= R(g, \Delta) \mid \tau \rightarrow \tau$$

guard-safe?      dependency on (transformed) samples

$$\frac{\Gamma \vdash L : R(t, \Delta) \quad \Gamma \vdash M : \sigma \quad \Gamma \vdash N : \sigma}{\Gamma \vdash \text{if } L < 0 \text{ then } M \text{ else } N : \sigma}$$

$$\tau ::= R(g, \Delta) \mid \tau \rightarrow \tau$$

guard-safe?      dependency on (transformed) samples

$$\frac{\Gamma \vdash L : R(t, \Delta) \quad \Gamma \vdash M : \sigma \quad \Gamma \vdash N : \sigma}{\Gamma \vdash \text{if } L < 0 \text{ then } M \text{ else } N : \sigma}$$

$$\overline{\Gamma \vdash 0 : R(f, \Delta)}$$

$$\tau ::= R(g, \Delta) \mid \tau \rightarrow \tau$$

guard-safe?      dependency on (transformed) samples

$$\frac{\Gamma \vdash L : R(\mathbf{t}, \Delta) \quad \Gamma \vdash M : \sigma \quad \Gamma \vdash N : \sigma}{\Gamma \vdash \text{if } L < 0 \text{ then } M \text{ else } N : \sigma}$$

$$\overline{\Gamma \vdash 0 : R(\mathbf{f}, \Delta)}$$

$$\overline{\Gamma \vdash \text{transform sample}_{\mathcal{N}} \text{ by } T : R(\mathbf{t}, \{s_j\})} \quad T \text{ diffeomorphic}$$

$$\tau ::= R(g, \Delta) \mid \tau \rightarrow \tau$$

guard-safe?  $\swarrow$   $\nwarrow$  dependency on (transformed) samples

$$\frac{\Gamma \vdash L : R(t, \Delta) \quad \Gamma \vdash M : \sigma \quad \Gamma \vdash N : \sigma}{\Gamma \vdash \text{if } L < 0 \text{ then } M \text{ else } N : \sigma}$$

$$\overline{\Gamma \vdash 0 : R(f, \Delta)}$$

$$\overline{\Gamma \vdash \text{transform sample}_{\mathcal{N}} \text{ by } T : R(t, \{s_j\})} \quad T \text{ diffeomorphic}$$

$\swarrow$  fresh



$\tau ::= R(\underline{g}, \Delta) \mid \tau \rightarrow \tau$   
 guard-safe?  $\swarrow$   $\nwarrow$  dependency on (transformed) samples

$$\frac{\Gamma \vdash L : R(\underline{t}, \Delta) \quad \Gamma \vdash M : \sigma \quad \Gamma \vdash N : \sigma}{\Gamma \vdash \text{if } L < 0 \text{ then } M \text{ else } N : \sigma}$$

$$\overline{\Gamma \vdash 0 : R(\underline{f}, \Delta)}$$

$\overline{\Gamma \vdash \text{transform sample}_{\mathcal{N}} \text{ by } T : R(\underline{t}, \{s_j\})}$   $T$  diffeomorphic  
 fresh  $\swarrow$

$$\frac{\Gamma \vdash M : R(\underline{t}, \Delta_1) \quad N : R(\underline{t}, \Delta_2)}{\Gamma \vdash M - N : R(\underline{t}, \Delta_1 \cup \Delta_2)} \quad \Delta_1 \cap \Delta_2 = \emptyset$$

$\tau ::= R(g, \Delta) \mid \tau \rightarrow \tau$   
 guard-safe?  $\nearrow$   $\nwarrow$  dependency on (transformed) samples

$$\frac{\Gamma \vdash L : R(\mathbf{t}, \Delta) \quad \Gamma \vdash M : \sigma \quad \Gamma \vdash N : \sigma}{\Gamma \vdash \text{if } L < 0 \text{ then } M \text{ else } N : \sigma}$$

$$\overline{\Gamma \vdash 0 : R(\mathbf{f}, \Delta)}$$

$$\overline{\Gamma \vdash \text{transform sample}_{\mathcal{N}} \text{ by } T : R(\mathbf{t}, \{s_j\})} \quad T \text{ diffeomorphic}$$

fresh  $\longleftarrow$

$$\frac{\Gamma \vdash M : R(\mathbf{t}, \Delta_1) \quad N : R(\mathbf{t}, \Delta_2)}{\Gamma \vdash M - N : R(\mathbf{t}, \Delta_1 \cup \Delta_2)} \quad \Delta_1 \cap \Delta_2 = \emptyset$$

*Establish correctness via logical relations*

## Uniform Convergence

If  $M$  is typable then

$$\mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})] \xrightarrow{\text{unif.}} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket(\boldsymbol{\theta}, \mathbf{s})] \quad \text{as } \eta \searrow 0 \text{ for } \boldsymbol{\theta} \in \Theta$$

## Uniform Convergence

If  $M$  is typable then

$$\mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})] \xrightarrow{\text{unif.}} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket(\boldsymbol{\theta}, \mathbf{s})] \quad \text{as } \eta \searrow 0 \text{ for } \boldsymbol{\theta} \in \Theta$$

For any error tolerance  $\epsilon > 0$ ,

exists accuracy coefficient  $\eta > 0$  s.t. **for all**  $\boldsymbol{\theta} \in \Theta$

$$\mathbb{E}_{\mathbf{s}}[\llbracket M \rrbracket(\boldsymbol{\theta}, \mathbf{s})] < \mathbb{E}_{\mathbf{s}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})] + \epsilon$$

## Uniform Convergence

If  $M$  is typable then

$$\mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})] \xrightarrow{\text{unif.}} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}}[\llbracket M \rrbracket(\boldsymbol{\theta}, \mathbf{s})] \quad \text{as } \eta \searrow 0 \text{ for } \boldsymbol{\theta} \in \Theta$$

For any error tolerance  $\epsilon > 0$ ,

exists accuracy coefficient  $\eta > 0$  s.t. **for all**  $\boldsymbol{\theta} \in \Theta$

$$\mathbb{E}_{\mathbf{s}}[\llbracket M \rrbracket(\boldsymbol{\theta}, \mathbf{s})] < \mathbb{E}_{\mathbf{s}}[\llbracket M \rrbracket_{\eta}(\boldsymbol{\theta}, \mathbf{s})] + \epsilon$$

In particular for  $\boldsymbol{\theta}^*$  obtained by

SGD with Reparameterisation Gradient (**fast!**) for  $\eta$ -smoothing

Part V:  
Empirical Evaluation

## Score Estimator

**X** high variance

## Standard Reparameterisation Estimator

**X** biased

## Score Estimator

**X** high variance

## Standard Reparameterisation Estimator

**X** biased

[Lee et al., NeurIPS 2018]:



## Score Estimator

✗ high variance

## Standard Reparameterisation Estimator

✗ biased

[Lee et al., NeurIPS 2018]:

- Fix bias with additional non-trivial *boundary* terms
- ✗ Only discuss efficient method for *affine* guards

## Score Estimator

✗ high variance

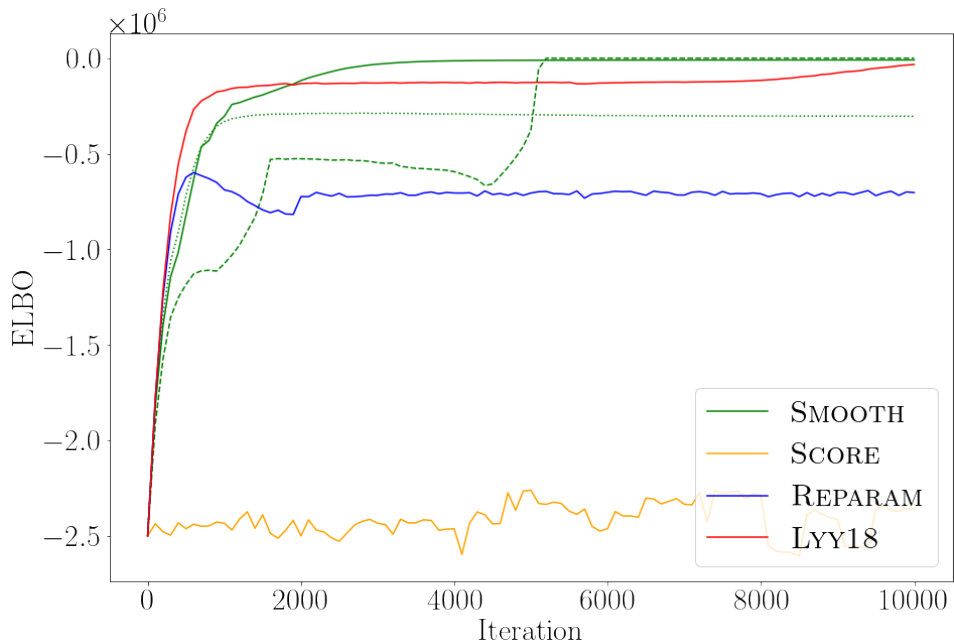
## Standard Reparameterisation Estimator

✗ biased

[Lee et al., NeurIPS 2018]:

- Fix bias with additional non-trivial *boundary* terms
- ✗ Only discuss efficient method for *affine* guards
- ✗ No discussion of PL aspects
- ✗ Only concerned with unbiasedness, not with overall *correctness* of SGD

temperature



## temperature: Variance and Cost

Estimator	Cost	Variance
Score	1	1
Reparam		
Smooth (ours)		
Lyy18		

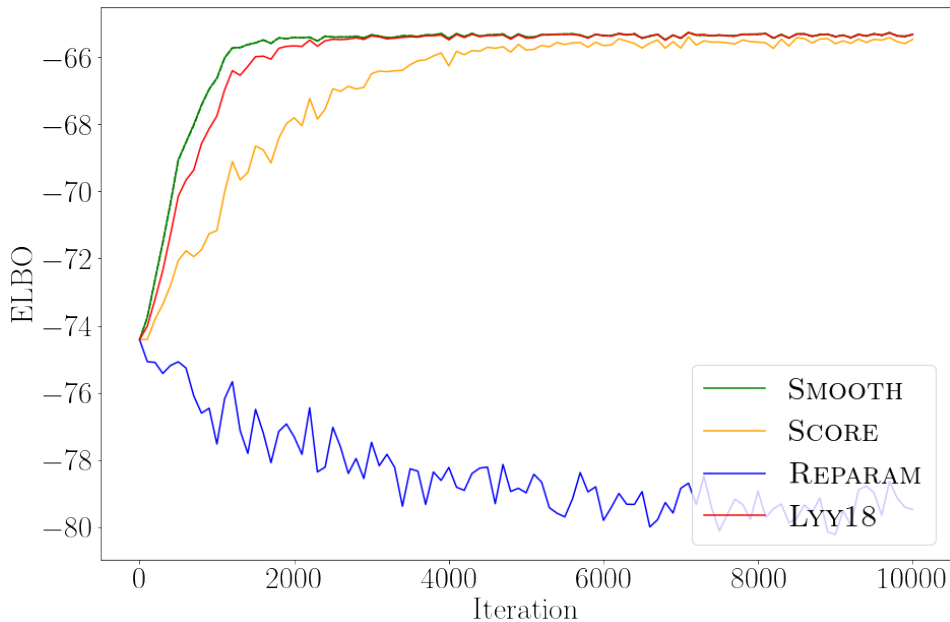
## temperature: Variance and Cost

Estimator	Cost	Variance
Score	1	1
Reparam	1.28	
Smooth (ours)	1.62	
Lyy18	<i>9.12</i>	

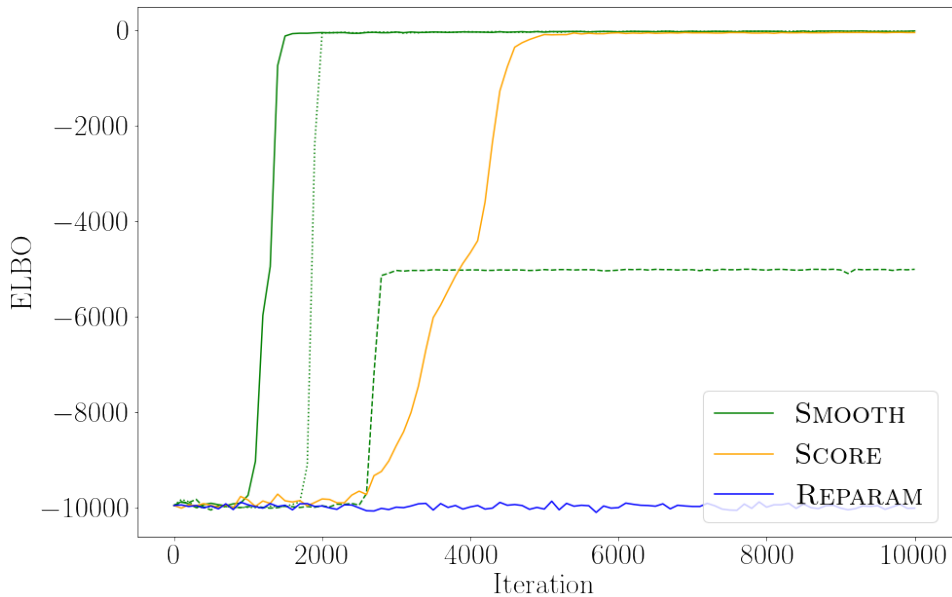
## temperature: Variance and Cost

Estimator	Cost	Variance
Score	1	1
Reparam	1.28	1.48e-08
Smooth (ours)	1.62	3.17e-10
Lyy18	<i>9.12</i>	1.22e-06

cheating



# xornet





# Fast and Correct Gradient-Based Optimisation for Probabilistic Programming via Smoothing

# Fast and Correct Gradient-Based Optimisation for Probabilistic Programming via Smoothing

- Smoothed semantics avoids bias (caused by branching)
  - ▶ categorical model based on Frölicher spaces

# Fast and Correct Gradient-Based Optimisation for Probabilistic Programming via Smoothing

- Smoothed semantics avoids bias (caused by branching)
  - ▶ categorical model based on Frölicher spaces
- Type systems enforce restrictions

# Fast and Correct Gradient-Based Optimisation for Probabilistic Programming via Smoothing

- Smoothed semantics avoids bias (caused by branching)
  - ▶ categorical model based on Frölicher spaces
- Type systems enforce restrictions
- Stochastic Gradient Descent is provably *correct*
- Approximations *converge* uniformly

# Fast and Correct Gradient-Based Optimisation for Probabilistic Programming via Smoothing

- Smoothed semantics avoids bias (caused by branching)
  - ▶ categorical model based on Frölicher spaces
- Type systems enforce restrictions
- Stochastic Gradient Descent is provably *correct*
- Approximations *converge* uniformly
- *Competitive* on benchmarks

# Fast and Correct Gradient-Based Optimisation for Probabilistic Programming via Smoothing

- Smoothed semantics avoids bias (caused by branching)
  - ▶ categorical model based on Frölicher spaces
- Type systems enforce restrictions
- Stochastic Gradient Descent is provably *correct*
- Approximations *converge* uniformly
- *Competitive* on benchmarks

## Ongoing Work

- Choice of accuracy coefficient