

Bayesian Inference

1. Exact Inference
2. Sampling-Based Methods: MCMC, HMC etc.
3. Variational Inference
 - ▶ frame posterior inference as (deterministic) *optimisation* problem

Variational Inference

Take: *variational family* $\{q_\theta \mid \theta \in \Theta\}$ of “simpler” *guide* distributions

Approach: find $\theta^* \in \Theta$ s.t. $q_\theta(\mathbf{z})$ is “closest” to $p(\mathbf{z} \mid \mathbf{x})$

KL divergence

or *maximise*

$$\text{ELBO}(\theta) := \mathbb{E}_{\mathbf{z} \sim q_\theta} [\log p(\mathbf{z}, \mathbf{x}) - \log q_\theta(\mathbf{z})]$$

model

guide

solve *optimisation* via **Stochastic Gradient Descent**

key ingredient!

Gradient Estimation (of Expectation)

- ▶ **Score Estimator:**
widely applicable but *high variance*
- ▶ **Reparameterisation Estimator:**
better in practice but *may be biased!* [Lee et al., NeurIPS 2018]

$$\nabla_\theta \mathbb{E}_{s \sim \mathcal{N}(0,1)} [\theta + s \geq 0] \neq \mathbb{E}_{s \sim \mathcal{N}(0,1)} [\underbrace{\nabla_\theta [\theta + s \geq 0]}_{=0 \text{ a.e.}}]$$

Contributions: Study reparameterisation gradient estimator for *continuous* but possibly *non-differentiable* programs

- ▶ categorical models
- ▶ prove *unbiasedness* in continuous setting
- ▶ establish *continuity* in languages with conditionals compositionally

Idealised Programming Language:

$M ::= x \mid r \mid f \mid \lambda x. M \mid MM$
 $\mid \text{sample } \mathcal{D}(M, \dots, M) \mid \text{observe } M \text{ from } \mathcal{D}(M, \dots, M)$
 $\mid \text{if } M < 0 \text{ then } M \text{ else } M \mid \text{sif } M < 0 \text{ then } M \text{ else } M$

smoothed conditional
[KOW23]

($r \in \mathbb{R}, f: \mathbb{R}^\ell \rightarrow \mathbb{R}$ and \mathcal{D} is a continuous probability distribution)

Example (Temperature Regulation)

- ▶ Without intervention, the temperature fluctuates randomly.
- ▶ If the temperature drops below a threshold of 18°C the heating is engaged and the power is proportional to the deviation from the threshold.
- ▶ Time is discretised and after one time unit we measure a temperature of 21°C.
- ▶ We are interested in the distribution of the original temperature.

```
let t0 = sample normal(20, sigma)
    mu = t0 + if t0 < 18 then c * (18 - t0)
              else 0
observe 21 from normal(mu, sigma)
in t0
```

($\sigma_0, \sigma, c > 0$ are constants.)

NB The joint density is not differentiable (yet continuous) at $t_0 = 18$.

Denotational Weight Semantics

Example. $\llbracket \text{sample } \mathcal{N}(20, \sigma_0) + x \cdot (\text{observe } 2 \text{ from } \mathcal{N}(0, 1)) \rrbracket (x, [s])$
 $= \text{pdf}_{\mathcal{N}}(s \mid 20, \sigma_0) \cdot \text{pdf}_{\mathcal{N}}(2 \mid 0, 1)$

- ▶ denotational version of weight semantics
- ▶ beyond measurability: capture piecewise definition and continuity
- ▶ complication: smoothed conditionals at higher-order [KOW23]

Solution: *Generalise Frölicher spaces*, replacing smoothness with functions $\mathbb{R} \rightarrow \mathbb{R}$ with mild closure properties, enriched over **Vect**.

1. Piecewise Analytic Functions under Analytic Partitions (PAP):

$$f(x) = \sum_{i=1}^{\ell} [x \in U_i] \cdot f_i(x)$$

($U_1, \dots, U_\ell \subseteq \mathbb{R}^n$ is a partition of analytic sets, each f_i is analytic) [LYRY20]

2. Continuous PAP (CPAP):

piecewise definitions agree on boundaries (for each $x \in \overline{U_i} \cap \overline{U_j}$, $f_i(x) = f_j(x)$)

Theorem (Unbiasedness)

If $f \circ \phi_\theta$ is a continuous PAP with partial derivatives which are uniformly dominated by an integrable function then

$$\nabla_\theta \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} [f(\phi_\theta(\mathbf{z}))] = \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} [\nabla_\theta f(\phi_\theta(\mathbf{z}))]$$

differentiable a.e.

Guarantee pre-conditions: densities are Schwartz functions, partial derivatives of ϕ_θ and primitives are bounded by polynomials

Interpret terms without conditionals (possibly with smoothed conditionals) in VectCPAP and obtain unbiasedness.

Example. Rephrase conditional via non-differentiable primitive:

$$c \cdot (\text{ReLU}(18 - t_0))$$

Continuity for Terms with Conditionals

$$\frac{\Gamma \vdash_{\text{cont}} L : R \quad \Gamma \vdash_{\text{cont}} M : \tau \quad \Gamma \vdash_{\text{cont}} N : \tau \quad \forall \gamma \in \llbracket \Gamma \rrbracket. \llbracket L \rrbracket(\gamma) = 0}{\Gamma \vdash_{\text{cont}} \text{if } L < 0 \text{ then } M \text{ else } N : \tau} \rightarrow \llbracket M \rrbracket(\gamma) = \llbracket N \rrbracket(\gamma)$$

Problem: *generally not tractable!*

Solution: *restrict guards to affine terms*

- ▶ efficiently sample from boundary
- ▶ efficiently check guard’s consistency (linear arithmetic solvers)

Example. For analytic primitives f and g ,

$$\text{if } x - y < 0 \text{ then } f \ x \ y \ \text{else } g \ x \ y \ \text{is continuous}$$

$$\iff (f - g)|_U = 0 \quad \text{where } U := \{(x, y) \mid x = y\}$$

$$\iff^* f(x, y) = g(x, y) \quad (x, y) \sim U$$

* with probability 1

Efficient Continuity Check

Suppose: $\text{if } (\underline{a}^T \mathbf{x} + \underline{c}) < 0 \text{ then } F \ \text{else } G$ w.l.o.g. $a_1 \neq 0$

Check: For all *consistent* branches S in F and T in G ,

$$\llbracket S \rrbracket(x_1, \dots, x_n) = \llbracket T \rrbracket(x_1, \dots, x_n)$$

where $x_2, \dots, x_n \sim \mathcal{D}$ and $x_1 := \frac{-a_2 x_2 + \dots + a_n x_n + c}{a_1}$.
(e.g. $\mathcal{N}(0, 1)$)

Example (Consistency)

Implement $\max\{|x|, 1\}$:

if $x + 1 < 0$ **then**

$-x$

else

(**if** $x - 1 < 0$ **then** 1 **else** x)

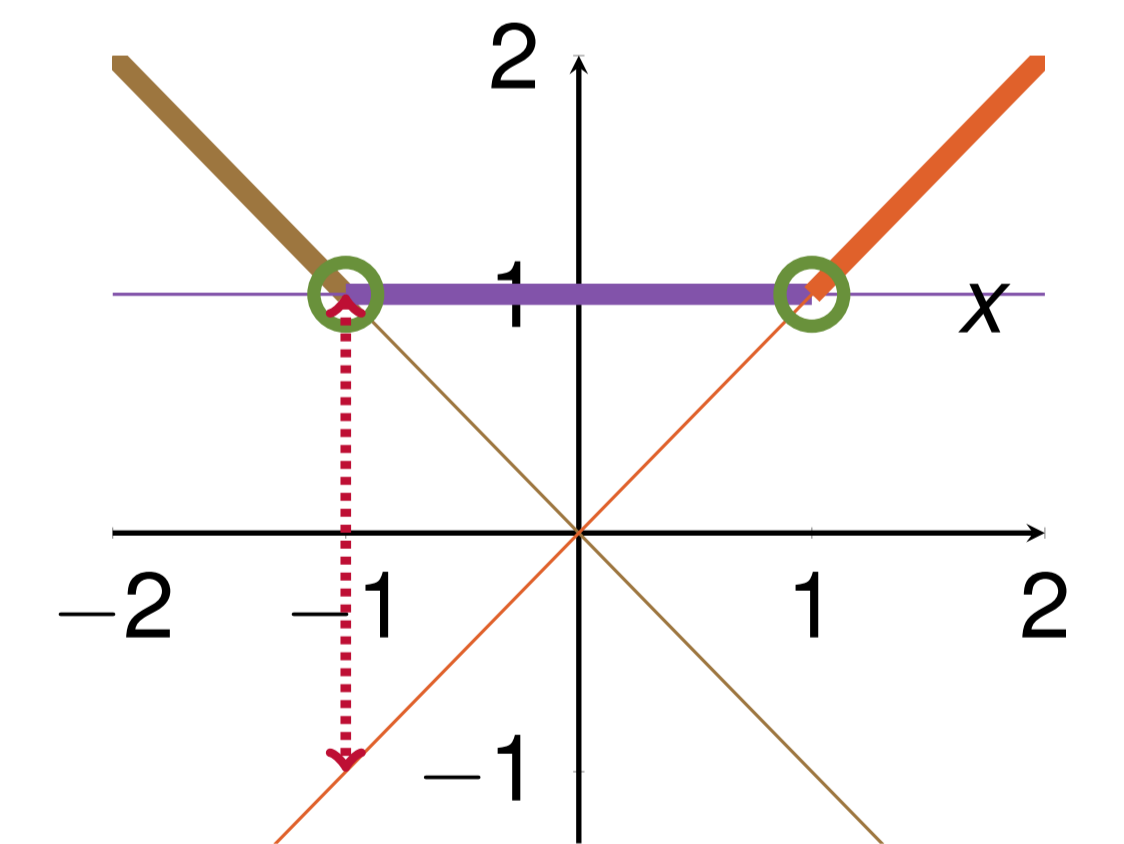
Sufficient to check:

$$\llbracket -x \rrbracket(-1) \stackrel{?}{=} \llbracket 1 \rrbracket(-1) \checkmark \quad (\text{outer conditional})$$

$$\llbracket 1 \rrbracket(1) \stackrel{?}{=} \llbracket x \rrbracket(1) \checkmark \quad (\text{inner conditional})$$

$(x + 1 \leq 0) \wedge (x - 1 \geq 0)$ is inconsistent, so *no* need to check:

$$\llbracket -x \rrbracket(-1) \stackrel{?}{=} \llbracket x \rrbracket(-1) \times \quad (\text{outer conditional})$$



References

[LYY18] Wonyeol Lee, Hangeol Yu, and Hongseok Yang: Reparameterization gradient for non-differentiable models. NeurIPS 2018.

[LYRY20] Wonyeol Lee, Hangeol Yu, Xavier Rival, and Hongseok Yang. On correctness of automatic differentiation for non-differentiable functions. NeurIPS 2020.

[KOW23] Basim Khajwal, C. H. Luke Ong, and Dominik Wagner. Fast and correct gradient-based optimisation for probabilistic programming via smoothing, ESOP 2023 to appear.