

Reparametrisation Tricks
or
How to Differentiate Through Samples from
Probability Distributions

Adam Kosiorek

May 24, 2017

Table of Contents

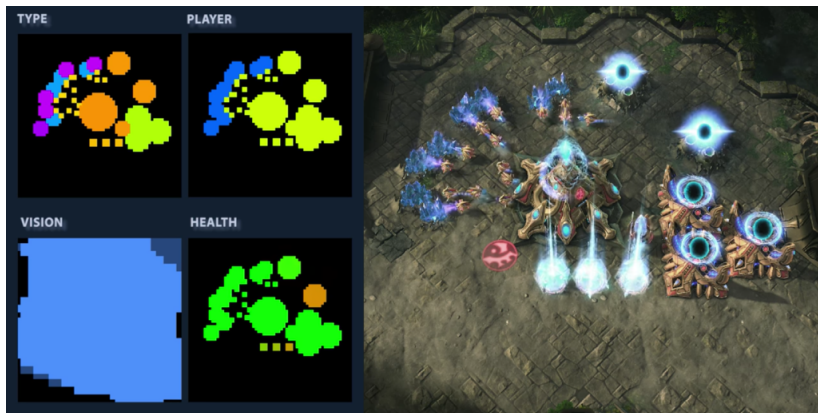
A Few Examples

Why do we need this?

Continuous Random Variables: One-Liners

Discrete Random Variables: REINFORCE

Playing Starcraft



$$\theta_t = f_{\phi}(\mathbf{x}_{1:t})$$

compute parameters

$$a_t \sim p(a \mid \theta_t)$$

sample action from a distribution

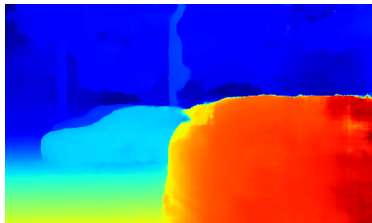
$$R = R(\mathbf{a}_{1:T}, \mathbf{x}_{1:T})$$

compute reward based on actions and states

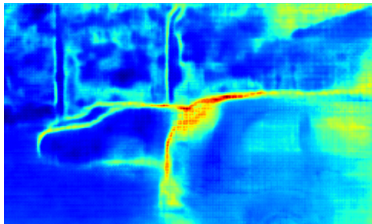
Alex Kendall's work on uncertainty



Input image I



Depth estimate d

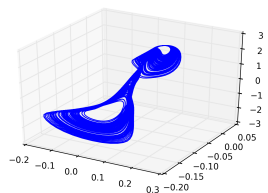
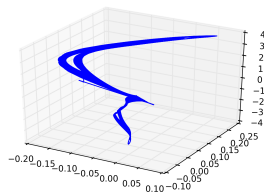
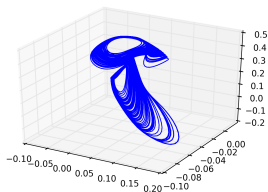


Uncertainty σ^2

$$\begin{aligned}\mu, \sigma^2 &= f_{\phi}(I) \\ d &\sim \mathcal{N}(d \mid \mu, \sigma^2)\end{aligned}$$

Stochastic Attractors of Chaotic Systems

Neil Dhir, Adam Kosiorek, Michael Osborne, Ingmar Posner



Dimensionality of the latent space?

$$d_E \sim q_\phi(d_E \mid \mathbf{x}_{1:T})$$

$$f_{d_E} : \mathbb{R}^N \rightarrow \mathbb{R}^{d_E}$$

$$f_{d_E} : \mathbf{x}_{1:T} \mapsto \mathbf{z}_{1:T}$$

Table of Contents

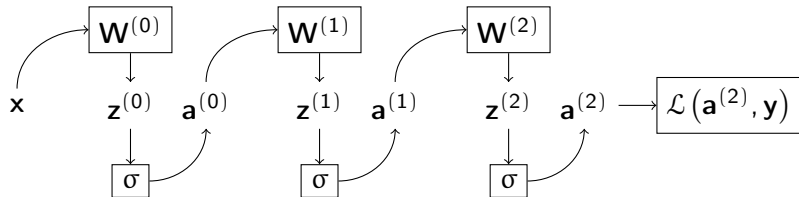
A Few Examples

Why do we need this?

Continuous Random Variables: One-Liners

Discrete Random Variables: REINFORCE

General Computation Graph



Backprop Equations

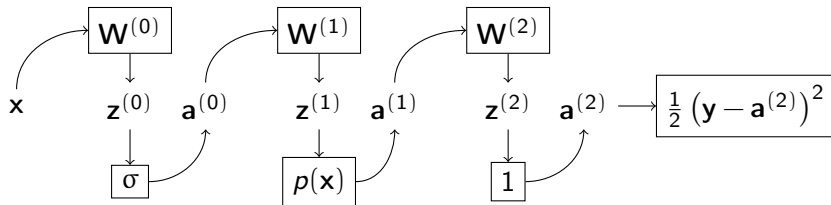
Let $\delta^{(l)} = \frac{\partial \mathcal{L}(\cdot)}{\partial \mathbf{z}^{(l)}}$, then backprop is given by:

$$\delta^{(L)} = \nabla_{\mathbf{a}^{(L)}} \mathcal{L}(\mathbf{a}^{(L)}, \mathbf{y}) \odot \sigma'(\mathbf{z}^{(L)}), \quad (1)$$

$$\delta^{(l)} = \left(\left(\mathbf{W}^{(l+1)} \right)^T \delta^{(l+1)} \right) \odot \sigma'(\mathbf{z}^{(l)}), \quad (2)$$

$$\frac{\partial \mathcal{L}(\cdot)}{\partial \mathbf{W}^{(l)}} = \delta^{(l)} \left(\mathbf{a}^{(l-1)} \right)^T. \quad (3)$$

Stochastic Computation Graph



Forward pass

$$\mathbf{z}^{(0)} = \mathbf{W}^{(0)} \mathbf{x}$$

Forward pass

$$\mathbf{z}^{(0)} = \mathbf{W}^{(0)} \mathbf{x}$$

$$\mathbf{a}^{(0)} = \sigma \left(\mathbf{z}^{(0)} \right)$$

Forward pass

$$\mathbf{z}^{(0)} = \mathbf{W}^{(0)} \mathbf{x}$$

$$\mathbf{a}^{(0)} = \sigma \left(\mathbf{z}^{(0)} \right)$$

$$\mathbf{z}^{(1)} = \mathbf{W}^{(1)} \mathbf{a}^{(0)}$$

Forward pass

$$\mathbf{z}^{(0)} = \mathbf{W}^{(0)} \mathbf{x}$$

$$\mathbf{a}^{(0)} = \sigma \left(\mathbf{z}^{(0)} \right)$$

$$\mathbf{z}^{(1)} = \mathbf{W}^{(1)} \mathbf{a}^{(0)}$$

$$\theta = \sigma \left(\mathbf{z}^{(1)} \right)$$

Forward pass

$$\mathbf{z}^{(0)} = \mathbf{W}^{(0)} \mathbf{x}$$

$$\mathbf{a}^{(0)} = \sigma \left(\mathbf{z}^{(0)} \right)$$

$$\mathbf{z}^{(1)} = \mathbf{W}^{(1)} \mathbf{a}^{(0)}$$

$$\theta = \sigma \left(\mathbf{z}^{(1)} \right)$$

$$\mathbf{a}^{(1)} \sim p(\mathbf{a} \mid \theta)$$

Forward pass

$$\mathbf{z}^{(0)} = \mathbf{W}^{(0)} \mathbf{x}$$

$$\mathbf{a}^{(0)} = \sigma \left(\mathbf{z}^{(0)} \right)$$

$$\mathbf{z}^{(1)} = \mathbf{W}^{(1)} \mathbf{a}^{(0)}$$

$$\theta = \sigma \left(\mathbf{z}^{(1)} \right)$$

$$\mathbf{a}^{(1)} \sim p(\mathbf{a} \mid \theta)$$

$$\mathbf{z}^{(2)} = \mathbf{W}^{(2)} \mathbf{a}^{(1)}$$

Forward pass

$$\mathbf{z}^{(0)} = \mathbf{W}^{(0)} \mathbf{x}$$

$$\mathbf{a}^{(0)} = \sigma \left(\mathbf{z}^{(0)} \right)$$

$$\mathbf{z}^{(1)} = \mathbf{W}^{(1)} \mathbf{a}^{(0)}$$

$$\theta = \sigma \left(\mathbf{z}^{(1)} \right)$$

$$\mathbf{a}^{(1)} \sim p(\mathbf{a} \mid \theta)$$

$$\mathbf{z}^{(2)} = \mathbf{W}^{(2)} \mathbf{a}^{(1)}$$

$$\mathbf{a}^{(2)} = \mathbf{z}^{(2)}$$

Forward pass

$$\mathbf{z}^{(0)} = \mathbf{W}^{(0)} \mathbf{x}$$

$$\mathbf{a}^{(0)} = \sigma \left(\mathbf{z}^{(0)} \right)$$

$$\mathbf{z}^{(1)} = \mathbf{W}^{(1)} \mathbf{a}^{(0)}$$

$$\theta = \sigma \left(\mathbf{z}^{(1)} \right)$$

$$\mathbf{a}^{(1)} \sim p(\mathbf{a} \mid \theta)$$

$$\mathbf{z}^{(2)} = \mathbf{W}^{(2)} \mathbf{a}^{(1)}$$

$$\mathbf{a}^{(2)} = \mathbf{z}^{(2)}$$

$$\mathcal{L}(\cdot) = \frac{1}{2} \left(\mathbf{y} - \mathbf{a}^{(2)} \right)^2$$

Backward pass

$$\mathbf{z}^{(0)} = \mathbf{W}^{(0)} \mathbf{x}$$

$$\mathbf{a}^{(0)} = \sigma \left(\mathbf{z}^{(0)} \right)$$

$$\mathbf{z}^{(1)} = \mathbf{W}^{(1)} \mathbf{a}^{(1)}$$

$$\theta = \sigma \left(\mathbf{z}^{(1)} \right)$$

$$\mathbf{a}^{(1)} \sim p(\mathbf{a} \mid \theta)$$

$$\mathbf{z}^{(2)} = \mathbf{W}^{(2)} \mathbf{a}^{(1)}$$

$$\mathbf{a}^{(2)} = \mathbf{z}^{(2)}$$

$$\delta^{(2)} = (\mathbf{y} - \mathbf{a}^{(2)})$$

$$\mathcal{L}(\cdot) = \frac{1}{2} \left(\mathbf{y} - \mathbf{a}^{(2)} \right)^2$$

Backward pass

$$\mathbf{z}^{(0)} = \mathbf{W}^{(0)} \mathbf{x}$$

$$\mathbf{a}^{(0)} = \sigma \left(\mathbf{z}^{(0)} \right)$$

$$\mathbf{z}^{(1)} = \mathbf{W}^{(1)} \mathbf{a}^{(1)}$$

$$\theta = \sigma \left(\mathbf{z}^{(1)} \right)$$

$$\mathbf{a}^{(1)} \sim p(\mathbf{a} \mid \theta)$$

$$\mathbf{z}^{(2)} = \mathbf{W}^{(2)} \mathbf{a}^{(1)}$$

$$\mathbf{a}^{(2)} = \mathbf{z}^{(2)}$$

$$\delta^{(2)} = (\mathbf{y} - \mathbf{a}^{(2)})$$

$$\mathcal{L}(\cdot) = \frac{1}{2} \left(\mathbf{y} - \mathbf{a}^{(2)} \right)^2$$

Backward pass

$$\mathbf{z}^{(0)} = \mathbf{W}^{(0)} \mathbf{x}$$

$$\mathbf{a}^{(0)} = \sigma(\mathbf{z}^{(0)})$$

$$\mathbf{z}^{(1)} = \mathbf{W}^{(1)} \mathbf{a}^{(1)} \quad \vdots$$

$$\theta = \sigma(\mathbf{z}^{(1)})$$

$$\mathbf{a}^{(1)} \sim p(\mathbf{a} \mid \theta) \quad \delta^{(1)} = \left(\left(\mathbf{W}^{(2)} \right)^T \delta^{(2)} \right) \odot \frac{\partial \mathbf{a}^{(1)}}{\partial \mathbf{z}^{(1)}}$$

$$\mathbf{z}^{(2)} = \mathbf{W}^{(2)} \mathbf{a}^{(1)}$$

$$\mathbf{a}^{(2)} = \mathbf{z}^{(2)} \quad \delta^{(2)} = (\mathbf{y} - \mathbf{a}^{(2)})$$

$$\mathcal{L}(\cdot) = \frac{1}{2} \left(\mathbf{y} - \mathbf{a}^{(2)} \right)^2$$

Gradient of a sample?

$$\begin{aligned}\frac{\partial \mathbf{a}^{(1)}}{\partial \mathbf{z}^{(1)}} &= \frac{\partial p(\mathbf{a}^{(1)} \mid \theta)}{\partial \mathbf{z}^{(1)}} \\ &= \frac{\partial p(\mathbf{a}^{(1)} \mid \theta)}{\partial \theta} \frac{\partial \theta}{\partial \mathbf{z}^{(1)}}\end{aligned}\tag{4}$$

Gradient of a sample?

$$\begin{aligned}\frac{\partial \mathbf{a}^{(1)}}{\partial \mathbf{z}^{(1)}} &= \frac{\partial p(\mathbf{a}^{(1)} \mid \theta)}{\partial \mathbf{z}^{(1)}} \\ &= \frac{\partial p(\mathbf{a}^{(1)} \mid \theta)}{\partial \theta} \frac{\partial \theta}{\partial \mathbf{z}^{(1)}}\end{aligned}\tag{4}$$

but...

$$\mathbf{a}^{(1)} \sim p(\mathbf{a} \mid \theta)\tag{5}$$

Gradient of a sample?

$$\begin{aligned}\frac{\partial \mathbf{a}^{(1)}}{\partial \mathbf{z}^{(1)}} &= \frac{\partial p(\mathbf{a}^{(1)} \mid \theta)}{\partial \mathbf{z}^{(1)}} \\ &= \frac{\partial p(\mathbf{a}^{(1)} \mid \theta)}{\partial \theta} \frac{\partial \theta}{\partial \mathbf{z}^{(1)}}\end{aligned}\tag{4}$$

but...

$$\begin{aligned}\mathbf{a}^{(1)} &\sim p(\mathbf{a} \mid \theta) \\ \frac{\partial p(\mathbf{a}^{(1)} \mid \theta)}{\partial \theta} &= \nexists\end{aligned}\tag{5}$$

Table of Contents

A Few Examples

Why do we need this?

Continuous Random Variables: One-Liners

Discrete Random Variables: REINFORCE

(Continuous) Reparametrisation Trick

Let $\mathbf{x} \sim p(\mathbf{x} \mid \theta)$ be a random variable. Perform change of variables such that

$$\begin{aligned}\epsilon &\sim p(\epsilon), \\ \mathbf{x} &= g(\epsilon, \theta)\end{aligned}\tag{6}$$

and treat ϵ as if it was a constant.

(Continuous) Reparametrisation Trick

Let $\mathbf{x} \sim p(\mathbf{x} \mid \theta)$ be a random variable. Perform change of variables such that

$$\begin{aligned}\epsilon &\sim p(\epsilon), \\ \mathbf{x} &= g(\epsilon, \theta)\end{aligned}\tag{6}$$

and treat ϵ as if it was a constant.

Example:

$$\begin{aligned}\epsilon &\sim \mathcal{N}(0, 1) \\ x &= 2 + \frac{1}{2}\epsilon \\ x &\sim \mathcal{N}\left(2, \frac{1}{4}\right) = \mathcal{N}(\mu, \sigma^2)\end{aligned}$$

(Continuous) Reparametrisation Trick

Let $\mathbf{x} \sim p(\mathbf{x} \mid \theta)$ be a random variable. Perform change of variables such that

$$\begin{aligned}\epsilon &\sim p(\epsilon), \\ \mathbf{x} &= g(\epsilon, \theta)\end{aligned}\tag{6}$$

and treat ϵ as if it was a constant.

Example:

$$\epsilon \sim \mathcal{N}(0, 1)$$

$$x = 2 + \frac{1}{2}\epsilon$$

$$x \sim \mathcal{N}\left(2, \frac{1}{4}\right) = \mathcal{N}(\mu, \sigma^2)$$

$$\frac{dx}{d\mu} = 1 \qquad \frac{dx}{d\sigma} = \frac{1}{2}\epsilon$$

Why can we do this?

We are NOT computing $\mathcal{L}(x)$.

Why can we do this?

We are NOT computing $\mathcal{L}(x)$.

We are computing $\mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] = \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \mathrm{d}x$.

Why can we do this?

We are computing $\mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] = \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \, dx$.

We get the gradient by the change of variables:

$$p(x) = \left| \frac{d\epsilon}{dx} \right| p(\epsilon) \quad \text{Change of variables,}$$

$$\implies |p(x) \, dx| = |p(\epsilon) \, d\epsilon| \quad \text{Mass conservation.}$$

Why can we do this?

We are computing $\mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] = \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \, dx$.

We get the gradient by the change of variables:

$$p(x) = \left| \frac{d\epsilon}{dx} \right| p(\epsilon) \quad \text{Change of variables,}$$

$$\implies |p(x) \, dx| = |p(\epsilon) \, d\epsilon| \quad \text{Mass conservation.}$$

$$\nabla_{\theta} \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] = \nabla_{\theta} \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \, dx$$

Why can we do this?

We are computing $\mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] = \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \, dx$.

We get the gradient by the change of variables:

$$\begin{aligned} p(x) &= \left| \frac{d\epsilon}{dx} \right| p(\epsilon) && \text{Change of variables,} \\ \implies |p(x) \, dx| &= |p(\epsilon) \, d\epsilon| && \text{Mass conservation.} \end{aligned}$$

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] &= \nabla_{\theta} \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \, dx \\ &= \nabla_{\theta} \int p(\epsilon)\mathcal{L}(x) \, d\epsilon \end{aligned}$$

Why can we do this?

We are computing $\mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] = \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \, dx$.

We get the gradient by the change of variables:

$$\begin{aligned} p(x) &= \left| \frac{d\epsilon}{dx} \right| p(\epsilon) && \text{Change of variables,} \\ \implies |p(x) \, dx| &= |p(\epsilon) \, d\epsilon| && \text{Mass conservation.} \end{aligned}$$

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] &= \nabla_{\theta} \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \, dx \\ &= \nabla_{\theta} \int p(\epsilon)\mathcal{L}(x) \, d\epsilon = \nabla_{\theta} \int p(\epsilon)\mathcal{L}(g(\epsilon, \theta)) \, d\epsilon \end{aligned}$$

Why can we do this?

We are computing $\mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] = \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \, dx$.

We get the gradient by the change of variables:

$$\begin{aligned} p(x) &= \left| \frac{d\epsilon}{dx} \right| p(\epsilon) && \text{Change of variables,} \\ \implies |p(x) \, dx| &= |p(\epsilon) \, d\epsilon| && \text{Mass conservation.} \end{aligned}$$

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] &= \nabla_{\theta} \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \, dx \\ &= \nabla_{\theta} \int p(\epsilon)\mathcal{L}(x) \, d\epsilon = \nabla_{\theta} \int p(\epsilon)\mathcal{L}(g(\epsilon, \theta)) \, d\epsilon \\ &= \int p(\epsilon) \nabla_{\theta} \mathcal{L}(g(\epsilon, \theta)) \, d\epsilon \end{aligned}$$

Why can we do this?

We are computing $\mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] = \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \, dx$.

We get the gradient by the change of variables:

$$\begin{aligned} p(x) &= \left| \frac{d\epsilon}{dx} \right| p(\epsilon) && \text{Change of variables,} \\ \implies |p(x) \, dx| &= |p(\epsilon) \, d\epsilon| && \text{Mass conservation.} \end{aligned}$$

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] &= \nabla_{\theta} \int p_{\theta}(x|\mathcal{D})\mathcal{L}(x) \, dx \\ &= \nabla_{\theta} \int p(\epsilon)\mathcal{L}(x) \, d\epsilon = \nabla_{\theta} \int p(\epsilon)\mathcal{L}(g(\epsilon, \theta)) \, d\epsilon \\ &= \int p(\epsilon) \nabla_{\theta} \mathcal{L}(g(\epsilon, \theta)) \, d\epsilon = \mathbb{E}_{p(\epsilon)}[\nabla_{\theta} \mathcal{L}(g(\epsilon, \theta))] \end{aligned}$$

MC approximation

How to compute $\mathbb{E}_{p(\epsilon)}[\nabla_{\theta} \mathcal{L}(g(\epsilon, \theta))]$?

1. Sample $\epsilon^{(s)} \sim p(\epsilon)$.
2. $\nabla_{\theta} \mathcal{L}(x) \simeq \frac{1}{S} \sum_{s=1}^S \nabla_{\theta} \mathcal{L}(g(\epsilon^{(s)}, \theta))$

Table of Contents

A Few Examples

Why do we need this?

Continuous Random Variables: One-Liners

Discrete Random Variables: REINFORCE

REINFORCE or the log-derivative trick

We can use REINFORCE for continuous variables, too, but one-liners are typically easier and implemented by default.

Recall that

$$\nabla \log f(x) = \frac{\nabla f(x)}{f(x)}.$$

REINFORCE

Again, we're computing the expectation of the gradient:

$$\nabla_{\theta} \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] = \int \nabla_{\theta} p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \mathrm{d}x$$

REINFORCE

Again, we're computing the expectation of the gradient:

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] &= \int \nabla_{\theta} p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \mathrm{d}x \\ &= \int \frac{p_{\theta}(x|\mathcal{D})}{p_{\theta}(x|\mathcal{D})} \nabla_{\theta} p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \mathrm{d}x\end{aligned}$$

REINFORCE

Again, we're computing the expectation of the gradient:

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] &= \int \nabla_{\theta} p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \, dx \\ &= \int \frac{p_{\theta}(x|\mathcal{D})}{p_{\theta}(x|\mathcal{D})} \nabla_{\theta} p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \, dx \\ &= \int p_{\theta}(x|\mathcal{D}) \nabla_{\theta} \log p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \, dx\end{aligned}$$

REINFORCE

Again, we're computing the expectation of the gradient:

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] &= \int \nabla_{\theta} p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \, dx \\&= \int \frac{p_{\theta}(x|\mathcal{D})}{p_{\theta}(x|\mathcal{D})} \nabla_{\theta} p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \, dx \\&= \int p_{\theta}(x|\mathcal{D}) \nabla_{\theta} \log p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \, dx \\&= \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\nabla_{\theta} \log p_{\theta}(x|\mathcal{D}) \mathcal{L}(x)]\end{aligned}$$

REINFORCE

Again, we're computing the expectation of the gradient:

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\mathcal{L}(x)] &= \int \nabla_{\theta} p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \, dx \\&= \int \frac{p_{\theta}(x|\mathcal{D})}{p_{\theta}(x|\mathcal{D})} \nabla_{\theta} p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \, dx \\&= \int p_{\theta}(x|\mathcal{D}) \nabla_{\theta} \log p_{\theta}(x|\mathcal{D}) \mathcal{L}(x) \, dx \\&= \mathbb{E}_{p_{\theta}(x|\mathcal{D})}[\nabla_{\theta} \log p_{\theta}(x|\mathcal{D}) \mathcal{L}(x)] \\&\simeq \frac{1}{S} \sum_{s=1}^S \nabla_{\theta} \log p_{\theta}(x|\mathcal{D}) \mathcal{L}(x)\end{aligned}$$