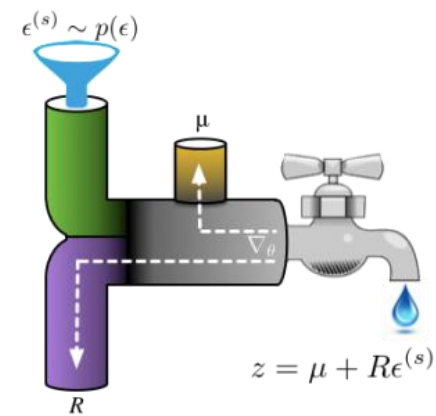


MC gradient estimators



Families of MC gradient estimators

- Score-function estimator
- Pathwise gradient estimators
- Measure-valued gradient estimators

Score-function estimator

- Better known as the REINFORCE algorithm
- Exploiting the following property

$$\frac{d}{d\mathbf{x}} \log f(\mathbf{x}) = \frac{1}{f(\mathbf{x})} \cdot \frac{df(\mathbf{x})}{d\mathbf{x}}$$

- When our function $f(\mathbf{x})$ is a probability density

$$\nabla_{\varphi} \log p_{\varphi}(\mathbf{x}) = \frac{1}{p_{\varphi}(\mathbf{x})} \nabla_{\varphi} p_{\varphi}(\mathbf{x}) \Leftrightarrow \nabla_{\varphi} p_{\varphi}(\mathbf{x}) = p_{\varphi}(\mathbf{x}) \nabla_{\varphi} \log p_{\varphi}(\mathbf{x})$$

- $\nabla_{\varphi} \log p(\mathbf{x})$: score-function
- A neat trick to rewrite the gradient of a density as another density

Deriving the score-function estimator

- As a use case the following expectation from VAE: $\nabla_{\phi} \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})]$

$$\begin{aligned} \nabla_{\phi} \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] &= \\ &= \nabla_{\phi} \int_{\mathbf{z}} \log p(\mathbf{x}|\mathbf{z}) q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\ &= \int_{\mathbf{z}} \log p(\mathbf{x}|\mathbf{z}) \nabla_{\phi} q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\ &= \int_{\mathbf{z}} \log p(\mathbf{x}|\mathbf{z}) q_{\phi}(\mathbf{z}|\mathbf{x}) \nabla_{\phi} \log q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\ &= \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z}) \nabla_{\phi} \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\ &\approx \frac{1}{n} \sum_i \log p(\mathbf{x}|\mathbf{z}^{(i)}) \nabla_{\phi} \log q_{\phi}(\mathbf{z}^{(i)}|\mathbf{x}), \mathbf{z}^{(i)} \sim q_{\phi}(\mathbf{z}|\mathbf{x}) \end{aligned}$$

Score-function estimator properties

- Any function $f(\mathbf{x})$ amenable
 - Good for simulators or black box functions (RL)
- The $p_\varphi(\mathbf{x})$ must be differentiable w.r.t. to parameters φ
- It must be easy to sample from $p_\varphi(\mathbf{x})$
- Unbiased estimator
- High variance estimator
 - The gradient will deviate a lot, but in the limit of many samples is accurate
 - Increases with more dimensions
 - If you sample once, this can be a problem and slow down or stop learning
 - Variance reduction methods are usually needed

Pathwise gradient estimator

- Also known as ‘reparameterization trick’
- Often the probability density can be rewritten as
 - a deterministic function of a simpler probability density
- Instead of sampling from a complex pdf \rightarrow sample from the simpler one
 - then transform deterministically the sample

$$\hat{\mathbf{x}} \sim p_{\varphi}(\mathbf{x}) \Leftrightarrow \hat{\mathbf{x}} = g(\hat{\boldsymbol{\varepsilon}}, \varphi), \hat{\boldsymbol{\varepsilon}} \sim p(\boldsymbol{\varepsilon})$$

Pathwise gradient estimator

$$\hat{\mathbf{x}} \sim p_{\varphi}(\mathbf{x}) \Leftrightarrow \hat{\mathbf{x}} = g(\hat{\boldsymbol{\varepsilon}}, \varphi), \hat{\boldsymbol{\varepsilon}} \sim p(\boldsymbol{\varepsilon})$$

- Stochasticity flows through a simple probability density
 - And, complexity flows from the deterministic transformation
 - For NN it means backprop –for deterministic functions only- is possible
- At the heart of this method is the change of variables formula

$$p_{\varphi}(\mathbf{x}) = p(\boldsymbol{\varepsilon}) |\det \nabla_{\boldsymbol{\varepsilon}} g(\boldsymbol{\varepsilon}, \varphi)|^{-1}$$

- We have seen normalizing flows using the same property

Deriving the pathwise gradient estimator

- As a use case the following expectation from VAE: $\nabla_{\varphi} \mathbb{E}_{\mathbf{z} \sim q_{\varphi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})]$
 - $\mathbf{z} = g(\boldsymbol{\varepsilon}, \varphi|\mathbf{x}) = \boldsymbol{\mu}_{\mathbf{x}} + \boldsymbol{\varepsilon} \cdot \boldsymbol{\sigma}_{\mathbf{x}}$, where $\varphi = (\boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\sigma}_{\mathbf{x}}) \Rightarrow d\mathbf{z} = \boldsymbol{\sigma}_{\mathbf{x}} d\boldsymbol{\varepsilon}$
 - $\det \nabla_{\boldsymbol{\varepsilon}} g(\boldsymbol{\varepsilon}, \varphi|\mathbf{x}) = \prod_i \sigma_{x,i}$

$$\begin{aligned} \nabla_{\varphi} \mathbb{E}_{\mathbf{z} \sim q_{\varphi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] &= \\ &= \nabla_{\varphi} \int_{\mathbf{z}} q_{\varphi}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}|\mathbf{z}) d\mathbf{z} \\ &= \nabla_{\varphi} \int_{\boldsymbol{\varepsilon}} \frac{1}{\prod_i \sigma_{x,i}} p(\boldsymbol{\varepsilon}) \log p(\mathbf{x}|g(\boldsymbol{\varepsilon}, \varphi|\mathbf{x})) \prod_i \sigma_{x,i} d\boldsymbol{\varepsilon} \\ &= \int_{\boldsymbol{\varepsilon}} p(\boldsymbol{\varepsilon}) \nabla_{\varphi} \log p(\mathbf{x}|g(\boldsymbol{\varepsilon}, \varphi|\mathbf{x})) d\boldsymbol{\varepsilon} \\ &= \mathbb{E}_{\boldsymbol{\varepsilon} \sim p(\boldsymbol{\varepsilon})} [\nabla_{\varphi} \log p(\mathbf{x}|g(\boldsymbol{\varepsilon}, \varphi|\mathbf{x}))] \\ &\approx \frac{1}{n} \sum_i \nabla_{\varphi} \log p(\mathbf{x}|g(\boldsymbol{\varepsilon}^{(i)}, \varphi|\mathbf{x})), \boldsymbol{\varepsilon}^{(i)} \sim p(\boldsymbol{\varepsilon}) \end{aligned}$$

Pathwise gradient estimator properties

- Only differentiable cost functions
 - Otherwise we cannot compute the $\nabla_{\varphi} f(\mathbf{x}, g(\boldsymbol{\varepsilon}, \varphi))$
 - Unlike score-function estimators that work with any cost function
- No need to know the pdf explicitly
 - Only the deterministic transformation and the base sampling distribution
- Low variance in general
 - Lower than the score-function estimator
 - Example: if you compare the VAE score-function and pathwise gradients, the score-function has an extra multiplicative term

$$\frac{1}{n} \sum_i \log p(\mathbf{x}|\mathbf{z}^{(i)}) \nabla_{\varphi} \log q_{\varphi}(\mathbf{z}^{(i)}|\mathbf{x}) \qquad \frac{1}{n} \sum_i \nabla_{\varphi} \log p(\mathbf{x}|g(\boldsymbol{\varepsilon}^{(i)}, \varphi))$$

- Very efficient (why proposed in VAE)
 - Even a single sample suffices no matter dimensionality

Qualitative comparison between estimators (1)

- Pathwise gradients have consistently lower variance

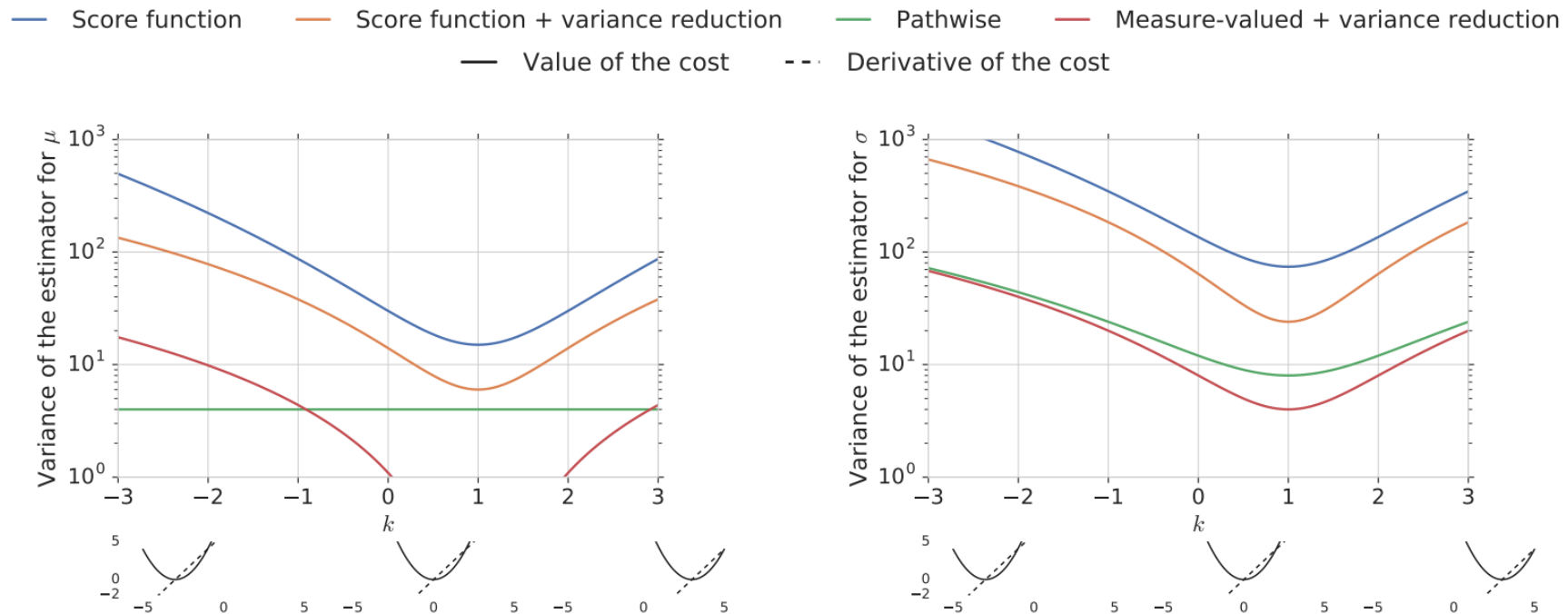


Figure 2: Variance of the stochastic estimates of $\nabla_{\theta} \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} [(x - k)^2]$ for $\mu = \sigma = 1$ as a function of k for three different classes of gradient estimators. Left: $\theta = \mu$; right: $\theta = \sigma$. The graphs in the bottom row show the function (solid) and its gradient (dashed) for $k \in \{-3, 0, 3\}$.

Qualitative comparison between estimators (2)

- For complex functions the pathwise gradient might have higher variance

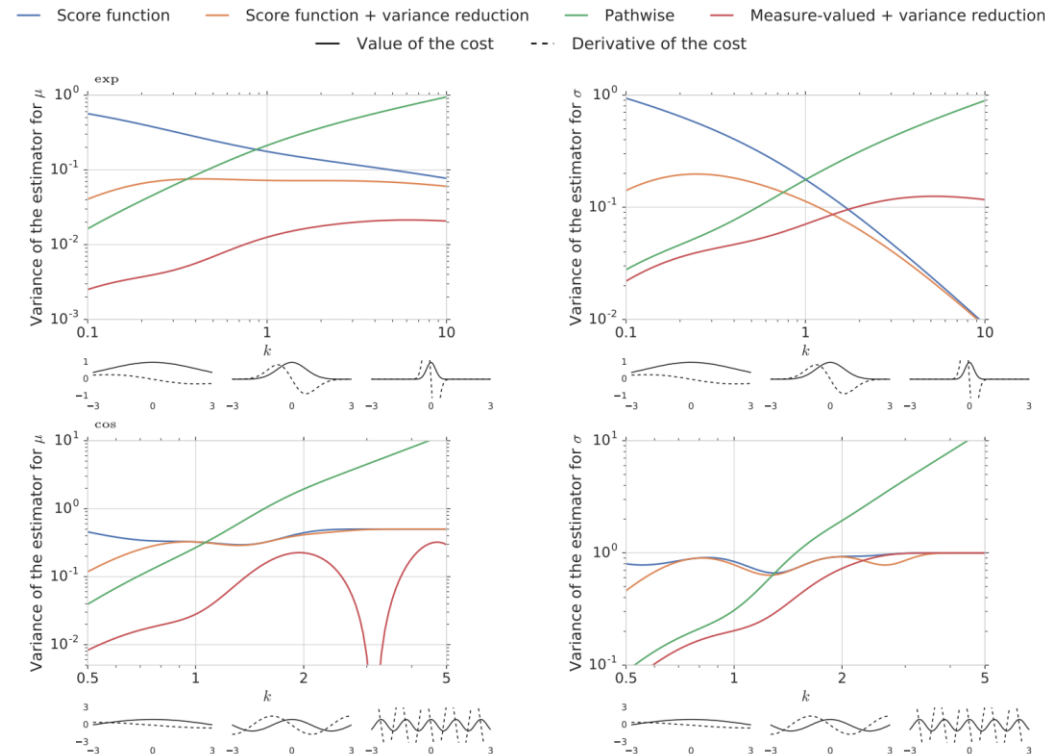


Figure 3: Variance of the stochastic estimates of $\nabla_{\theta} \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} [f(x; k)]$ for $\mu = \sigma = 1$ as a function of k . Top: $f(x; k) = \exp(-kx^2)$, bottom: $f(x; k) = \cos kx$. Left: $\theta = \mu$; right: $\theta = \sigma$. The graphs in the bottom row show the function (solid) and its gradient (dashed): for $k \in \{0.1, 1, 10\}$ for the exponential function, and $k \in \{0.5, 1.58, 5\}$ for the cosine function.

[Monte Carlo Gradient Estimation in Machine Learning](#)