

Approximating Gradients of Expectations and It's Applications

Deniz Koyuncu

December 4, 2019

- First Half: Approximating Gradients of Expectations
 - A Review, Approximating Expectations
 - A General Objective Function
 - VI Objective Function
 - Score Functions
 - Pathwise Gradient Estimators
- Second Half: Application to VAE
 - Models with Unobserved Variables
 - EM Algorithm
 - VAE Model Description
 - Connection to EM

Approximating Expectations

$$E[X] = \int P(X = x)xdx$$

- How can we empirically calculate the expectation?
- Assume X_1, \dots, X_N are i.i.d samples of X .

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$$

$$E[\hat{\mu}] = \frac{1}{N} \sum_{i=1}^N E[X_i] = E[X] \quad (1)$$

$$\text{Var}(\hat{\mu}) = \frac{1}{N^2} \sum_{i=1}^N \text{Var}(X_i) = \frac{\text{Var}(X)}{N}$$

- $\hat{\mu}$ is unbiased and notice that for finite variance $\lim_{N \rightarrow \infty} \text{Var}(\hat{\mu}) = 0$

Expectation of a Function of a R.V

$$E[g(X)] = \int P(X = x)g(x)dx$$

- Similarly, assume X_1, \dots, X_N are i.i.d samples of X .

$$\hat{\mu}_g = \frac{1}{N} \sum_{i=1}^N g(X_i)$$

$$E[\hat{\mu}_g] = \frac{1}{N} \sum_{i=1}^N E[g(X_i)] = E[g(X)] \quad (2)$$

$$\text{Var}(\hat{\mu}_g) = \frac{1}{N^2} \sum_{i=1}^N \text{Var}(g(X_i)) = \frac{\text{Var}(g(X))}{N}$$

- Again, $\hat{\mu}_g$ is unbiased and for finite variance $\lim_{N \rightarrow \infty} \text{Var}(\hat{\mu}_g) = 0$

A Specific Form of Objective Function

- Following [1]'s derivation,
- In our discussion \mathbf{X} is assumed to be continuous random vector
- Today, we will be talking about the objective functions in the following form

$$\begin{aligned} E_{P(\mathbf{X};\theta)}[g(\mathbf{X})] &= \int P(\mathbf{X} = \mathbf{x}; \theta)g(\mathbf{x})d\mathbf{x} \\ \eta &= \nabla_{\theta}E_{P(\mathbf{X};\theta)}[g(\mathbf{X})] = \nabla_{\theta} \int P(\mathbf{X} = \mathbf{x}; \theta)g(\mathbf{x})d\mathbf{x} \\ &= \int \nabla_{\theta}P(\mathbf{X} = \mathbf{x}; \theta)g(\mathbf{x})d\mathbf{x} \end{aligned} \tag{3}$$

- Notice that Eq. 3, is no longer an expectation and the integral may not have a closed form
- We're interested in somehow turning it into an expectation so that we can empirically estimate it

VI Objective Function

- VI tries to approximate $P(\mathbf{X}_h | \mathbf{X}_o = \mathbf{x}_o)$, with another distribution when exact calculation is not an option.
- One distance measure is the following KL divergence:

$$q^* = \arg \min_q KL(q(\mathbf{X}_h) || P(\mathbf{X}_h | \mathbf{X}_o = \mathbf{x}_o)) \quad (4)$$

- It's minimum when they're equal.
- The generally intractable term $P(\mathbf{X}_o = \mathbf{x}_o)$ is a constant for q .

Gradient Based VI

- Previously we were optimizing in the function space q ,
- In order to use gradient based optimization, we parameterize q as $q(\mathbf{X}_h; \boldsymbol{\lambda})$ where $\boldsymbol{\lambda}$ are the parameters,
- Therefore, the optimization function becomes,

$$\boldsymbol{\lambda}^* = \arg \min_{\boldsymbol{\lambda}} KL(q(\mathbf{X}_h; \boldsymbol{\lambda}) \parallel P(\mathbf{X}_h \mid \mathbf{X}_o = \mathbf{x}_o)) \quad (5)$$

- One can take the gradient as follows

Calculating the Gradients

$$\begin{aligned} & \nabla_{\lambda} KL(q(\mathbf{X}_h; \lambda) \parallel P(\mathbf{X}_h \mid \mathbf{X}_o = \mathbf{x}_o)) \\ &= \nabla_{\lambda} \int q(\mathbf{X}_h = \mathbf{x}_h; \lambda) \log \frac{q(\mathbf{X}_h = \mathbf{x}_h; \lambda)}{P(\mathbf{X}_h = \mathbf{x}_h, \mathbf{X}_o = \mathbf{x}_o)} d\mathbf{x}_h \\ &= -\nabla_{\lambda} \int q(\mathbf{X}_h = \mathbf{x}_h; \lambda) \log P(\mathbf{X}_h = \mathbf{x}_h, \mathbf{X}_o = \mathbf{x}_o) d\mathbf{x}_h \\ &+ \nabla_{\lambda} \int q(\mathbf{X}_h = \mathbf{x}_h; \lambda) \log q(\mathbf{X}_h = \mathbf{x}_h; \lambda) d\mathbf{x}_h \end{aligned} \tag{6}$$

- Observe that the second term is the - entropy of q ,

$$H_q(\mathbf{X}; \lambda) = - \int q(\mathbf{X}_h = \mathbf{x}_h; \lambda) \log q(\mathbf{X}_h = \mathbf{x}_h; \lambda) d\mathbf{x}_h$$

- If q is a exponential family than entropy has an analytical form.

$$\begin{aligned}\nabla_{\lambda} \int q(\mathbf{X}_h = \mathbf{x}_h; \lambda) \log P(\mathbf{X}_h = \mathbf{x}_h, \mathbf{X}_o = \mathbf{x}_o) d\mathbf{x}_h \\ &= \nabla_{\lambda} \int q(\mathbf{X}_h = \mathbf{x}_h; \lambda) g(\mathbf{x}_h) d\mathbf{x}_h \\ &= \nabla_{\lambda} \mathbb{E}_{q(\mathbf{x}_h; \lambda)} [g(\mathbf{X}_h)]\end{aligned}\tag{7}$$

- It's in the same form with η .

2 Approaches

- We're going to be talking about two approaches to approximate the gradient η ,
 - 1 Score Functions
 - 2 Pathwise Gradient Estimators

- One way of turning Eq. 3 into an expectation is as follows

$$\begin{aligned}\nabla_{\theta} E_{P(\mathbf{X}; \theta)}[g(\mathbf{X})] &= \int \nabla_{\theta} P(\mathbf{X} = \mathbf{x}; \theta) g(\mathbf{x}) d\mathbf{x} \\ &= \int \frac{P(\mathbf{X} = \mathbf{x}; \theta)}{P(\mathbf{X} = \mathbf{x}; \theta)} \nabla_{\theta} P(\mathbf{X} = \mathbf{x}; \theta) g(\mathbf{x}) d\mathbf{x} \quad (8) \\ &= E_{P(\mathbf{X}; \theta)} \left[\frac{\nabla_{\theta} P(\mathbf{X}; \theta)}{P(\mathbf{X}; \theta)} g(\mathbf{X}) \right]\end{aligned}$$

- Now, since it's a expectation we can do the empirical estimation.

- Using $\nabla \log g(\mathbf{x}) = \frac{\nabla g(\mathbf{x})}{g(\mathbf{x})}$

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{P(\mathbf{X}; \theta)}[g(\mathbf{X})] &= \mathbb{E}_{P(\mathbf{X}; \theta)}[\nabla_{\theta} \log P(\mathbf{X}; \theta) g(\mathbf{X})] \\ &= \mathbb{E}_{P(\mathbf{X}; \theta)}[\vec{h}(\mathbf{X})]\end{aligned}\tag{9}$$

$$\begin{aligned}\hat{\eta} &= \frac{1}{N} \sum_{s=1}^N \vec{h}(\mathbf{x}^{(s)}) \\ &= \frac{1}{N} \sum_{s=1}^N \nabla_{\theta} \log P(\mathbf{X} = \mathbf{x}^{(s)}; \theta) g(\mathbf{x}^{(s)}),\end{aligned}\tag{10}$$

where $\mathbf{x}^{(s)} \sim P(\mathbf{X}; \theta)$

- $\hat{\eta}$ is an unbiased estimator of η .
- As N increases, variance of gradient of each parameter decrease.
- There are interesting interpretations of the variance of score function estimator one can refer to [1].
- There are different approaches to reduce the variance of the score function estimator.

Pathwise Gradient Estimators

- In cases where the random vector \mathbf{X} can be written as a function of another random vector ϵ ,

$$\mathbf{X} = t(\epsilon; \theta)$$

- And one can sample from ϵ , we can use pathwise gradient estimators.
- One example is MVN [1],

$$\begin{aligned}\mathbf{X} &\sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ \epsilon &\sim N(\mathbf{0}, \mathbf{I}) \\ \mathbf{X} = t(\epsilon; \theta) &= \boldsymbol{\mu} + \mathbf{L}\epsilon, \quad \text{where } \mathbf{L}\mathbf{L}^T = \boldsymbol{\Sigma}\end{aligned}\tag{11}$$

- The following approach can be used.

$$\begin{aligned} E_{P(\mathbf{X};\theta)}[g(\mathbf{X})] &= E_{P(\epsilon)}[g(t(\epsilon; \theta))] \\ &= \int P(\epsilon = \varepsilon)g(t(\varepsilon; \theta))d\varepsilon \end{aligned} \tag{12}$$

- Notice that the distribution the expectation is with respect to namely $P(\epsilon)$, is independent of θ this time.

- Taking the gradient,

$$\begin{aligned}
 \nabla_{\theta} \mathbb{E}_{P(\mathbf{x}; \theta)}[g(\mathbf{X})] &= \nabla_{\theta} \int P(\epsilon = \varepsilon) g(t(\varepsilon; \theta)) d\varepsilon \\
 &= \int P(\epsilon = \varepsilon) \nabla_{\theta} g(t(\varepsilon; \theta)) d\varepsilon \\
 &= \mathbb{E}_{P(\epsilon)}[\vec{p}(\epsilon)]
 \end{aligned} \tag{13}$$

$$\begin{aligned}
 \hat{\eta} &= \frac{1}{N} \sum_{s=1}^N \vec{p}(\varepsilon^{(s)}) \\
 &= \frac{1}{N} \sum_{s=1}^N g(t(\varepsilon^{(s)}; \theta))
 \end{aligned} \tag{14}$$

where $\varepsilon^{(s)} \sim P(\epsilon)$

- $\hat{\eta}$ is an unbiased estimator of η and the variance goes to zero for each parameter.
- A more interesting remark is the comparison of two gradient estimators we have covered.
 - In [1], one remark made is the bound on the variance of the pathwise gradient estimator doesn't depend on the number of parameters.
 - The bound of variance of score function estimator depends on the number of parameters.
 - This doesn't imply that pathwise approach always have lower variance.
[1]
- A meaningful question to ask can be which distributions can be written in terms other simpler distributions.



Shakir Mohamed et al. *Monte Carlo Gradient Estimation in Machine Learning*. 2019. arXiv: 1906.10652 [stat.ML].