

(cont'd)

score

$$\begin{aligned}
 & \nabla_{\psi} E_{s \sim q} f(s) \\
 &= \nabla_{\psi} \sum_s q(s) f(s) \\
 &= \sum_s \nabla_{\psi} [q(s) f(s)] \\
 &= \sum_s \left( [\nabla_{\psi} q(s)] f(s) + q(s) \nabla_{\psi} f(s) \right) \\
 &= \sum_s [q(s) \nabla \log q(s)] f(s) + \cancel{q(s) \nabla f(s)} \\
 &= \sum_s q(s) [f(s) \nabla \log q(s)] + \cancel{\nabla f(s)} \\
 &= E_{s \sim q} [f(s) \cdot \nabla \log q(s)] + \cancel{\nabla f(s)} + \cancel{\nabla f(s)}
 \end{aligned}$$

~~score~~

Note:  $\nabla q(s) \cdot \nabla \log q(s) = \frac{\nabla q(s)}{q(s)}$

$$\approx \frac{1}{L} \sum_{l=1}^L [f(s^l) \nabla \log q(s^l) + \cancel{\nabla f(s^l)}]$$

where  $s^l \sim q$

~~$\frac{1}{L} \sum_{l=1}^L \nabla f(s^l)$  is fine.~~

but  $\frac{1}{L} \sum_{l=1}^L f(s^l) \nabla \log q(s^l)$

has high variance,

i.e., its value varies widely depending on the samples  $s^l$ .

why?

# Example

(2)

why?

suppose  $\mathcal{X} = \mathcal{V}$ ,  $s \in \{0, 1\}$

$$q(s=1) = \alpha$$

$$q(s=0) = 1 - \alpha$$

$$\therefore \nabla_{\alpha}^{\log} q(s=1) = \frac{1}{\alpha} > 0$$

$$\nabla_{\alpha}^{\log} q(s=0) = \frac{-1}{1-\alpha} < 0$$

if  $L=1$ , then

$$f(s) \nabla \log q(s)$$

has 2 values

$$f(1) / \alpha$$

$$-f(0) / (1-\alpha)$$

(22)

if  $f(1) + f(0)$  are both large  
positive numbers, then  $f(1)/2 \gg 0$   
&  $-f(0)/(1-\alpha) \ll 0$

ii, the value of  $f(s) \nabla \log q(s)$   
varies widely, even if  $f(s)$   
does not.

# More Generally

---

$$\frac{1}{L} \sum_{s^x=1}^L f(s^x) \nabla \log q(s^x)$$

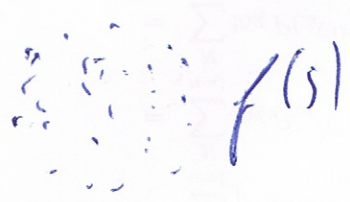
↑ discrete vector.

has high variance because  $f(s^x)$  is large (& ~~negative~~) for all  $s^x$ , because it is a likelihood, whereas  $\nabla \log q(s^x)$  is <sup>log</sup> +ve for some  $s^x$  & -ve for others, since  $q$  is a distribution, so an ~~increase~~ increase in ~~some~~  $q$  for some  $s^x$  must result in a decrease for other  $s^x$ .

So,  $f(s^x) \cdot \nabla \log q(s^x)$  is large & pos. for some  $s^x$ , & large & neg for other  $s^x$ .  
ie, it has high variance.

In practice

222



$f(s) \cdot \nabla \log q(s)$

56

high variance

# Centering

23

We avoid this by centering  $f$ .

i.e., replace  $f(s)$  by  $f(s) - b$   
for some constant,  $b$ .

note:  $\nabla \mathbb{E}_{s \sim q} [f(s) - b]$

$$= \nabla \left( \mathbb{E}_{s \sim q} f(s) - b \right)$$

$$= \nabla \mathbb{E}_{s \sim q} f(s)$$

Choose  $b$  so that  $f(s) - b$  does not have large values of

~~$f(1) \cdot \frac{f(1) - b}{2} = \frac{f(1) - b}{1 - \alpha}$~~  same sign  $\nabla$ s

~~$(f(1) - b) \nabla \log q(1)$   
 $= (f(0) - b) \nabla \log q(0)$~~

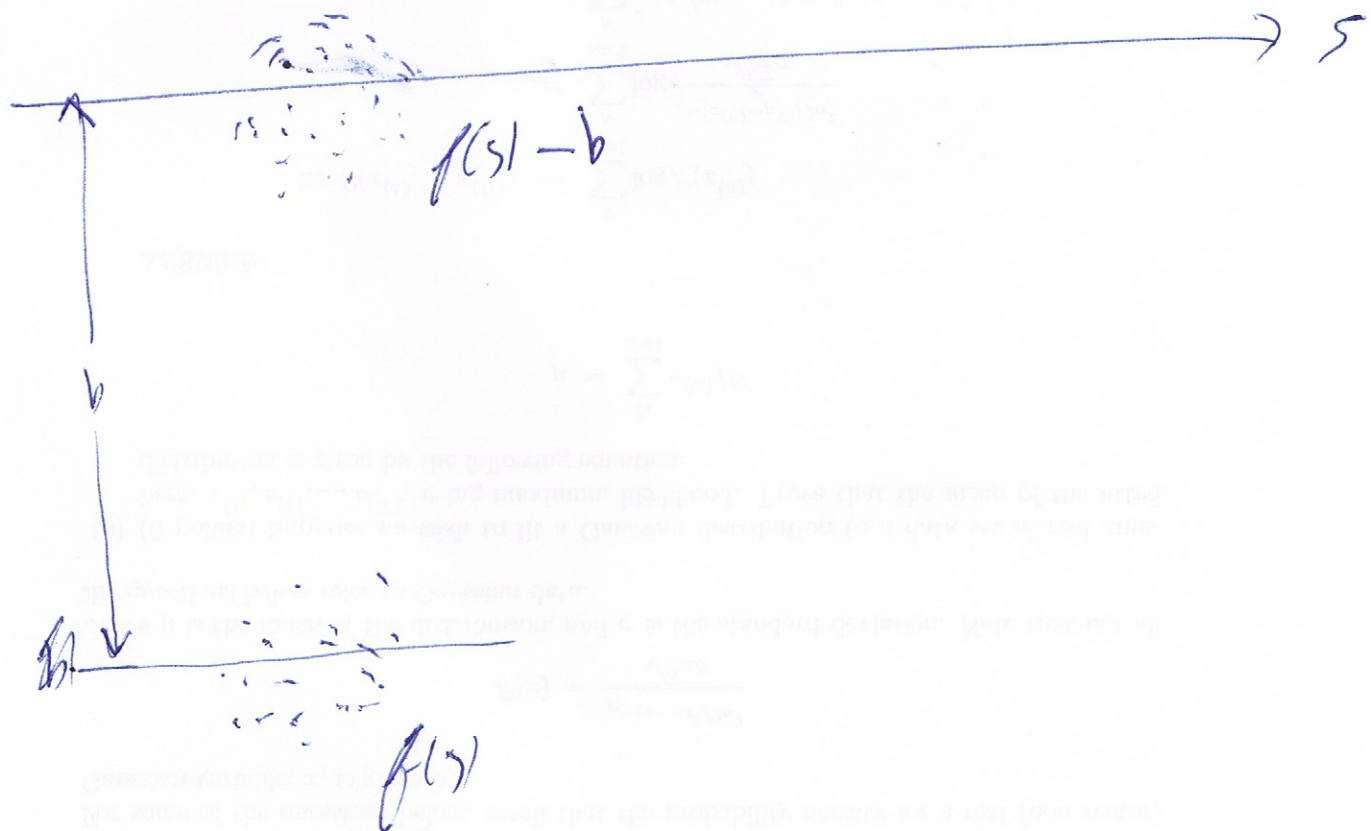
~~exercise exercise~~



We choose  $b$  so that  $f(s) - b$  has values ~~at~~ near 0. This gives

$(f(s) - b) \nabla \log q(s)$  low variance.

eg. choose  $b = \frac{1}{L} \sum_{i=1}^L f(s^i)$ .



(But what if  $L=1$ ?)

In our case,

This is an example of a  
Variance reduction technique.

Statistics has many.

(eg. Control variates.)

# Neural Variational Inference (25)

---

Recall:

$$\mathcal{L} = \sum_n E_{s \sim q(s|x^n)} f(s, x^n)$$

$$\nabla \mathcal{L} = \sum_n \nabla_{\theta} E_{s \sim q(s|x^n)} f(s, x^n)$$

~~of  $f$  &  $q$~~   
depend on  $\theta$

$$= \sum_n E_{s \sim q(s|x^n)} [f(s, x^n) \cdot \nabla_{\theta} q(s|x^n)]$$

~~\* easy stuff~~

$$= \sum_n E_{s \sim q} [f(s, x^n) - b_n] \cdot \nabla q(s|x^n)$$

$b_n$  depends on  $x^n$ , just since  $f + q$  do, i.e.,  $b_n = b(x^n)$

Use a neural net with input  $x$  to learn a good value of  $b$ . (see paper).

$b$  is called a baseline

of  $\nabla L$

The variance  $\checkmark$  is high if  $f(s, x^n)$  has a large range tends to be large & of the same sign for most values of  $s$ .

So, NUI tries to shrink the average value of  $[f(s, x^n) - b(x^n)]^2$ , averaged over  $s$ .

ie, to make  $b(x^n) \approx E_{s \sim q^n} f(s, x^n)$

Note:  $f(s) \cdot \nabla \log q(s)$

can be computed even if  $f$  is non-differentiable.

if, we can estimate

$$\nabla_{\psi} E_{s \sim q} f(s) \cdot \log q(s)$$

even without having to backpropagate through  $f$ .