

# Score functions

Xun Zheng  
xzheng1@andrew.cmu.edu

## 1 Score functions

### 1.1 Fisher score function

- Random variable  $x$  has density  $f_\theta(x)$  parameterized by  $\theta$ . Assume the support  $\mathcal{X}$  is independent of  $\theta$ .
- Score function measures the sensitivity of log-likelihood  $\log f_\theta(x)$  to its parameter  $\theta$ :

$$s(\hat{\theta}, x) := \frac{\partial \log f_{\hat{\theta}}(x)}{\partial \theta} = \frac{1}{f_{\hat{\theta}}(x)} \frac{\partial f_{\hat{\theta}}(x)}{\partial \theta} \quad (1)$$

Note that it is a function of both  $\theta$  and  $x$ .

- Mean of score is zero:

$$\mathbb{E}_{x \sim f_{\hat{\theta}}}[s(\hat{\theta}, x)] = \int_{\mathcal{X}} \frac{\partial}{\partial \theta} f_{\hat{\theta}}(x) dx = \frac{\partial}{\partial \theta} \underbrace{\int_{\mathcal{X}} f_{\hat{\theta}}(x) dx}_{=1} = 0 \quad (2)$$

Since the domain  $\mathcal{X}$  is independent of  $\theta$ , by Leibniz rule the integral and differentiation can be switched. Note that  $\theta$  for expectation and score function has to be the same.

- Variance of score is the Fisher information:

$$\mathcal{I}(\hat{\theta}) := \text{Var}_{x \sim f_{\hat{\theta}}}[s(\hat{\theta}, x)] = -\mathbb{E}_{x \sim f_{\hat{\theta}}}\left[\frac{\partial^2}{\partial \theta^2} \log f_{\hat{\theta}}(x)\right] \quad (3)$$

which is also negative average Hessian of log density. To see why,

$$\text{LHS} = \mathbb{E}_{x \sim f_{\hat{\theta}}}[s(\hat{\theta}, x)^2] - \underbrace{(\mathbb{E}_{x \sim f_{\hat{\theta}}}[s(\hat{\theta}, x)])^2}_{=0} \quad (4)$$

$$\text{RHS} = \mathbb{E}_{x \sim f_{\hat{\theta}}}\left[\underbrace{\left(\frac{1}{f_{\hat{\theta}}(x)} \frac{\partial f_{\hat{\theta}}(x)}{\partial \theta}\right)^2}_{=s(\hat{\theta}, x)}\right] - \underbrace{\mathbb{E}_{x \sim f_{\hat{\theta}}}\left[\frac{\partial^2 f_{\hat{\theta}}(x)}{\partial \theta^2} \frac{1}{f_{\hat{\theta}}(x)}\right]}_{=0 \text{ by Leibniz}} \quad (5)$$

- Informal summary:

$$\begin{array}{ccccc} \log f & \xrightarrow{\partial_\theta} & \overset{s}{(\log f)'} & \xrightarrow{\partial_\theta} & (\log f)'' \\ & & \downarrow \mathbb{E}_x & \searrow \text{Var}_x & \downarrow -\mathbb{E}_x \\ & & 0 & & \mathcal{I} \end{array}$$

## 1.2 Hyvärinen score function

- Hyvärinen score measures the sensitivity of log-likelihood w.r.t.  $x$ :

$$h_{f_\theta}(\hat{x}) := \frac{\partial \log f_\theta(\hat{x})}{\partial x} = \frac{1}{f_\theta(\hat{x})} \frac{\partial f_\theta(\hat{x})}{\partial x} \quad (6)$$

- Hyvärinen (2005) phrased this as the Fisher score w.r.t. a hypothetical location parameter: define a new density  $f_{\theta, \mu}(x) := f_\theta(x - \mu)$ , and let  $\hat{\mu} = 0$ , then the Fisher score of  $f_{\theta, \mu}$  restricted to  $\mu$  is

$$s(\hat{\mu}, \hat{x}) = \frac{\partial \log f_\theta(\hat{x} - \hat{\mu})}{\partial \mu} = \frac{1}{f_\theta(\hat{x} - \hat{\mu})} \frac{\partial f_\theta(\hat{x} - \hat{\mu})}{\partial(x - \mu)} \frac{\partial(\hat{x} - \hat{\mu})}{\partial \mu} = -h_{f_\theta}(\hat{x}) \quad (7)$$

Caution: for some distributions  $\mu$  may affect the support  $\mathcal{X}$ , such as exponential distribution, in which case the regularity condition for (2) does not hold anymore:

$$f_\theta(x) = \theta e^{-\theta x} \cdot \mathbb{I}(x > 0) \implies \mathbb{E}_{x \sim f_\theta} [h_{f_\theta}(x)] = -\theta \neq 0 \quad (8)$$

- Hyvärinen score is agnostic about the partition function. Suppose

$$f_\theta(x) = \frac{1}{Z_\theta} \tilde{f}_\theta(x), \quad Z_\theta = \int_{\mathcal{X}} \tilde{f}_\theta(x) dx \quad (9)$$

then the Hyvärinen score is just the sensitivity of the log-unnormalized-density w.r.t.  $x$ :

$$h_{f_\theta}(\hat{x}) = \frac{\partial [\log \tilde{f}_\theta(\hat{x}) - \log Z_\theta]}{\partial x} = \frac{\partial \log \tilde{f}_\theta(\hat{x})}{\partial x} \quad (10)$$

- Score matching can be useful in the following parametric density estimate problem: given i.i.d. samples from some unknown density  $f^*$ , estimate a parametric density  $f_\theta$ , only using evaluations of  $\tilde{f}_\theta$ .
- Score matching density estimator:  $\hat{\theta} = \operatorname{argmin}_\theta J(\theta)$ , where

$$J(\theta) = \mathbb{E}_{x \sim f^*} \left[ \frac{1}{2} \|h_{f_\theta}(x) - h_{f^*}(x)\|_2^2 \right] \quad (11)$$

$$= \mathbb{E}_{x \sim f^*} \left[ \frac{1}{2} \|h_{f_\theta}(x)\|_2^2 + \underbrace{\frac{1}{2} \|h_{f^*}(x)\|_2^2}_{\text{independent of } \theta} - \underbrace{\langle h_{f_\theta}(x), h_{f^*}(x) \rangle}_\alpha \right] \quad (12)$$

$$= \mathbb{E}_{x \sim f^*} \left[ \operatorname{tr}(H) + \frac{1}{2} \|\mathbf{h}\|^2 \right] + \text{constant} \quad (13)$$

with shortcuts

$$\mathbf{h} = h_{f_\theta}(x), \quad H = \frac{\partial h_{f_\theta}(x)}{\partial x} = \frac{\partial^2 \log f_\theta(x)}{\partial x \partial x} \quad (14)$$

$$\log f \xrightarrow{\partial_x} \underbrace{(\log f)'}_{\mathbf{h}} \xrightarrow{\partial_x} \underbrace{(\log f)''}_{H}$$

To see how  $f^*$  disappears, (roughly)

$$\mathbb{E}[\alpha] = \int f^*(x) \frac{\partial \log f_\theta(x)}{\partial x} \frac{\partial \log f^*(x)}{\partial x} dx \quad (15)$$

$$= \int f^*(x) \frac{\partial \log f_\theta(x)}{\partial x} \left( \frac{1}{f^*(x)} \frac{\partial f^*(x)}{\partial x} \right) dx \quad (16)$$

$$= \int \frac{\partial \log f_\theta(x)}{\partial x} \frac{\partial f^*(x)}{\partial x} dx \quad (17)$$

$$= - \int \frac{\partial^2 \log f_\theta(x)}{\partial x \partial x} f^*(x) dx \quad (18)$$

Last line follows from integration by parts and diminishing tails:  $\int g f' = [g f]_{-\infty}^{+\infty} - \int g' f$ .

Under regularity conditions,  $J(\theta) = 0 \iff \theta = \theta^*$ . Since  $J(\theta) \geq 0$ , this motivates the score matching density estimator.

- Pros:  $J(\theta)$  can be optimized using only gradient and Hessian of  $\log \tilde{f}_\theta(x)$ , no estimates of  $Z_\theta$  and  $f^*$  are involved.
- Cons: need to know how  $\theta$  behaves in  $\mathbf{h}$  and  $\text{tr}(H)$ , *e.g.*, access to  $\frac{\partial \text{tr}(H)}{\partial \theta}$ ,  $\frac{\partial \mathbf{h}}{\partial \theta}$ .

## References

- A. Hyvärinen. Estimation of Non-Normalized Statistical Models by Score Matching. *Journal of Machine Learning Research*, 2005.