Supp. Material: Reward Augmented Maximum Likelihood for Neural Structured Prediction

Zhifeng Chen Mohammad Norouzi Samy Bengio Navdeep Jaitly **Mike Schuster** Yonghui Wu **Dale Schuurmans** {mnorouzi, bengio, zhifengc, ndjaitly}@google.com {schuster, yonghui, schuurmans}@google.com

Google Brain

A Proofs

for some a

Proposition 1. For any twice differentiable strictly convex closed potential F, and $p, q \in int(\mathcal{F})$:

$$D_F(q \parallel p) = D_F(p \parallel q) + \frac{1}{4}(p-q)^{\mathsf{T}} \big(H_F(a) - H_F(b) \big)(p-q)$$
(1)
= $(1-\alpha)p + \alpha q, \ (0 < \alpha < \frac{1}{2}), \ b = (1-\beta)q + \beta p, \ (0 < \beta < \frac{1}{2}).$

Proof. Let f(p) denote $\nabla F(p)$ and consider the midpoint $\frac{q+p}{2}$. One can express $F(\frac{q+p}{2})$ by two Taylor expansions around p and q. By Taylor's theorem there is an $a = (1-\alpha)p + \alpha q$ for $0 \le \alpha \le \frac{1}{2}$ and $b = \beta p + (1 - \beta)q$ for $0 \le \beta \le \frac{1}{2}$ such that

$$F(\frac{q+p}{2}) = F(p) + (\frac{q+p}{2} - p)^{\top} f(p) + \frac{1}{2} (\frac{q+p}{2} - p)^{\top} H_F(a) (\frac{q+p}{2} - p)$$
(2)

$$= F(q) + \left(\frac{q+p}{2} - q\right)^{\top} f(q) + \frac{1}{2} \left(\frac{q+p}{2} - q\right)^{\top} H_F(b) \left(\frac{q+p}{2} - q\right), \quad (3)$$

$$2F(\frac{q+p}{2}) = 2F(p) + \left(q - p\right)^{\top} f(p) + \frac{1}{4} \left(q - p\right)^{\top} H_F(a) \left(q - p\right) \quad (4)$$

$$= 2F(q) + (p-q)^{\top}f(q) + \frac{1}{4}(p-q)^{\top}H_F(b)(p-q).$$
(5)

(4)

 \square

Therefore,

$$F(p) + F(q) - 2F(\frac{q+p}{2}) = F(p) - F(q) - (p-q)^{\top} f(q) - \frac{1}{4} (p-q)^{\top} H_F(b)(p-q)$$
(6)
$$= F(q) - F(q) - (q-q)^{\top} f(q) - \frac{1}{4} (p-q)^{\top} H_F(b)(p-q)$$
(7)

$$= F(q) - F(p) - (q-p)'f(p) - \frac{1}{4}(q-p)'H_F(a)(q-p)$$
(7)

$$= D_F(p \parallel q) - \frac{1}{4}(p-q)^{\top} H_F(b)(p-q)$$
(8)

$$= D_F(q \parallel p) - \frac{1}{4}(q-p)^{\top} H_F(a)(q-p),$$
(9)

leading to the result.

For the proof of Proposition 2, we first need to introduce a few definitions and background results. A Bregman divergence is defined from a strictly convex, differentiable, closed potential function F: $\mathcal{F} \to \mathbb{R}$, whose strictly convex conjugate $F^* : \mathcal{F}^* \to \mathbb{R}$ is given by $F^*(r) = \sup_{r \in \mathcal{F}} \langle r, q \rangle - F(q)$ [1]. Each of these potential functions have corresponding transfers, $f : \mathcal{F} \to \mathcal{F}^*$ and $f^* : \mathcal{F}^* \to \mathcal{F}$, given by the respective gradient maps $f = \nabla F$ and $f^* = \nabla F^*$. A key property is that $f^* = f^{-1}$ [1], which allows one to associate each object $q \in \mathcal{F}$ with its transferred image $r = f(q) \in \mathcal{F}^*$ and vice versa. The main property of Bregman divergences we exploit is that a divergence between any two domain objects can always be equivalently expressed as a divergence between their transferred images; that is, for any $p \in \mathcal{F}$ and $q \in \mathcal{F}$, one has [1]:

$$D_F(p || q) = F(p) - \langle p, r \rangle + F^*(r) = D_{F^*}(r || s), \qquad (10)$$

$$D_F(q \| p) = F^*(s) - \langle s, q \rangle + F(q) = D_{F^*}(s \| r), \qquad (11)$$

where s = f(p) and r = f(q). These relations also hold if we instead chose $s \in \mathcal{F}^*$ and $r \in \mathcal{F}^*$ in the range space, and used $p = f^*(s)$ and $q = f^*(r)$. In general (10) and (11) are not equal.

Two special cases of the potential functions F and F^* are interesting as they give rise to KL divergences. These two cases include $F_{\tau}(p) = -\tau \mathbb{H}(p)$ and $F_{\tau}^*(s) = \tau lse(s/\tau) =$ $\tau \log \sum_{y} \exp(s(y)/\tau)$, where $lse(\cdot)$ denotes the log-sum-exp operator. The respective gradient maps are $f_{\tau}(p) = \tau(\log(p) + 1)$ and $f_{\tau}^*(s) = f^*(s/\tau) = \frac{1}{\sum_y \exp(s(y)/\tau)} \exp(s/\tau)$, where f_{τ}^* denotes the normalized exponential operator for $\frac{1}{\tau}$ -scaled logits. Below, we derive $D_{F_{\tau}^*}(r \parallel s)$ for such F^*_{τ} :

$$D_{F_{\tau}^{*}}(s \parallel r) = F_{\tau}^{*}(s) - F_{\tau}^{*}(r) - (s - r)^{\mathsf{T}} \nabla F_{\tau}^{*}(r)$$

$$= \tau \operatorname{lse}(s/\tau) - \tau \operatorname{lse}(r/\tau) - (s - r)^{\mathsf{T}} f_{\tau}^{*}(r)$$

$$= -\tau \left((s/\tau - \operatorname{lse}(s/\tau)) - (r/\tau - \operatorname{lse}(r/\tau)) \right)^{\mathsf{T}} f_{\tau}^{*}(r)$$

$$= \tau f_{\tau}^{*}(r)^{\mathsf{T}} \left((r/\tau - \operatorname{lse}(r/\tau)) - (s/\tau - \operatorname{lse}(s/\tau)) \right)$$

$$= \tau f_{\tau}^{*}(r)^{\mathsf{T}} \left(\log f_{\tau}^{*}(r) - \log f_{\tau}^{*}(s) \right)$$

$$= \tau D_{\mathrm{KL}} (f_{\tau}^{*}(r) \parallel f_{\tau}^{*}(s))$$

$$= \tau D_{\mathrm{KL}} (q \parallel p)$$
(12)

Proposition 2. The KL divergence between p and q in two directions can be expressed as, $D_{\mathrm{KL}}(p \parallel q) = D_{\mathrm{KL}}(q \parallel p) + \frac{1}{4\tau^2} \operatorname{Var}_{y \sim f^*(a/\tau)} [s(y) - r(y)] - \frac{1}{4\tau^2} \operatorname{Var}_{y \sim f^*(b/\tau)} [s(y) - r(y)] = \frac{1}{4\tau^2} \operatorname{Var}_{y \sim f^*(b/\tau)} [s(y) - r(y)]$ $< D_{\mathrm{KL}}(q \parallel p) + \frac{1}{\tau^2} \|s - r\|_2^2,$ (14)

for some $a = (1 - \alpha)s + \alpha r$, $(0 \le \alpha \le \frac{1}{2})$, $b = (1 - \beta)r + \beta s$, $(0 \le \beta \le \frac{1}{2})$.

Proof. First, for the potential function $F_{\tau}^*(r) = \tau lse(r/\tau)$ it is easy to verify that F_{τ}^* satisfies the conditions for Proposition 1, and

$$H_{F^*_{\tau}}(a) = \frac{1}{\tau} (\text{Diag}(f^*_{\tau}(a)) - f^*_{\tau}(a) f^*_{\tau}(a)^{\top}) , \qquad (15)$$

where Diag(v) returns a square matrix the main diagonal of which comprises a vector v. Therefore, by Proposition 1 we obtain

$$D_{F_{\tau}^{*}}(r \parallel s) = D_{F_{\tau}^{*}}(s \parallel r) + \frac{1}{4}(s-r)^{\top}(H_{F_{\tau}^{*}}(a) - H_{F_{\tau}^{*}}(b))(s-r), \qquad (16)$$

for some $a = (1 - \alpha)s + \alpha r$, $(0 \le \alpha \le \frac{1}{2})$, $b = (1 - \beta)r + \beta s$, $(0 \le \beta \le \frac{1}{2})$. Note that by the specific form (15) we also have

$$(s-r)^{\top} H_{F_{\tau}^{*}}(a)(s-r) = \frac{1}{\tau} (s-r)^{\top} \left(\text{Diag}(f_{\tau}^{*}(a)) - f_{\tau}^{*}(a) f_{\tau}^{*}(a)^{\top} \right) (s-r)$$
(17)
$$= \frac{1}{\tau} \left(E_{\mathbf{v}_{\tau}, f^{*}(a)} \left[(s(\mathbf{v}) - r(\mathbf{v}))^{2} \right] - E_{\mathbf{v}_{\tau}, f^{*}(a)} \left[s(\mathbf{v}) - r(\mathbf{v}) \right]^{2} \right)$$
(18)

$$= \frac{1}{\tau} \left(E_{\mathbf{y} \sim f_{\tau}^*(a)} \left[(s(\mathbf{y}) - r(\mathbf{y}))^2 \right] - E_{\mathbf{y} \sim f_{\tau}^*(a)} \left[s(\mathbf{y}) - r(\mathbf{y}) \right]^2 \right)$$
(18)

$$= \frac{1}{\tau} \operatorname{Var}_{\mathbf{y} \sim f_{\tau}^{*}(a)} \left[s(\mathbf{y}) - r(\mathbf{y}) \right] , \qquad (19)$$

and
$$(s-r)^{\top} H_{F_{\tau}^*}(b)(s-r) = \frac{1}{\tau} \operatorname{Var}_{\mathbf{y} \sim f_{\tau}^*(b)} [s(\mathbf{y}) - r(\mathbf{y})]$$
. (20)
Therefore, by combining (19) and (20) with (16) we obtain

$$D_{F_{\tau}^{*}}(r \parallel s) = D_{F_{\tau}^{*}}(s \parallel r) + \frac{1}{4\tau} \operatorname{Var}_{\mathbf{y} \sim f_{\tau}^{*}(a)}[s(\mathbf{y}) - r(\mathbf{y})] - \frac{1}{4\tau} \operatorname{Var}_{\mathbf{y} \sim f_{\tau}^{*}(b)}[s(\mathbf{y}) - r(\mathbf{y})] .$$
(21)

Equality (13) then follows by applying (12) to (21).

Next, to prove the inequality in (14), let $\delta = s - r$ and observe that

$$D_{F_{\tau}^{*}}(r \parallel s) - D_{F_{\tau}^{*}}(s \parallel r) = \frac{1}{4}\delta^{\top} \left(H_{F_{\tau}^{*}}(a) - H_{F_{\tau}^{*}}(b) \right)\delta$$
(22)

$$= \frac{1}{4\tau} \delta^{\top} \text{Diag}(f_{\tau}^{*}(a) - f_{\tau}^{*}(b))\delta + \frac{1}{4\tau} \left(\delta^{\top} f_{\tau}^{*}(b)\right)^{2} - \frac{1}{4\tau} \left(\delta^{\top} f_{\tau}^{*}(a)\right)^{2}$$
(23)

$$\leq \frac{1}{4\tau} \|\delta\|_2^2 \|f_\tau^*(a) - f_\tau^*(b)\|_\infty + \frac{1}{4\tau} \|\delta\|_2^2 \|f_\tau^*(b)\|_2^2 \tag{24}$$

$$\leq \frac{1}{2\tau} \|\delta\|_2^2 + \frac{1}{4\tau} \|\delta\|_2^2 \tag{25}$$

since $||f_{\tau}^{*}(a) - f_{\tau}^{*}(b)||_{\infty} \leq 2$ and $||f_{\tau}^{*}(b)||_{2}^{2} \leq ||f_{\tau}^{*}(b)||_{1}^{2} \leq 1$. The result follows by applying (12) to (25).

References

[1] A. Banerjee, S. Merugu, I. S. Dhillon, and J. Ghosh. Clustering with Bregman divergences. JMLR, 2005.