

---

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

---

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

## A Proofs

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

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

for some  $a = (1 - \alpha)p + \alpha q$ , ( $0 \leq \alpha \leq \frac{1}{2}$ ),  $b = (1 - \beta)q + \beta p$ , ( $0 \leq \beta \leq \frac{1}{2}$ ).

*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  $0 \leq \alpha \leq \frac{1}{2}$  and  $b = \beta p + (1 - \beta)q$  for  $0 \leq \beta \leq \frac{1}{2}$  such that

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

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

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

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

Therefore,

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

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

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

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

leading to the result.  $\square$

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  $F : \mathcal{F} \rightarrow \mathbb{R}$ , whose strictly convex conjugate  $F^* : \mathcal{F}^* \rightarrow \mathbb{R}$  is given by  $F^*(r) = \sup_{q \in \mathcal{F}} \langle r, q \rangle - F(q)$  [1]. Each of these potential functions have corresponding transfers,  $f : \mathcal{F} \rightarrow \mathcal{F}^*$  and  $f^* : \mathcal{F}^* \rightarrow \mathcal{F}$ , 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]:

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

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

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  $F_\tau^*(s) = \tau \text{lse}(s/\tau) = \tau \log \sum_y \exp(s(y)/\tau)$ , where  $\text{lse}(\cdot)$  denotes the log-sum-exp operator. The respective gradient maps are  $f_\tau(p) = \tau(\log(p) + \mathbf{1})$  and  $f_\tau^*(s) = f^*(s/\tau) = \frac{1}{\sum_y \exp(s(y)/\tau)} \exp(s/\tau)$ , where  $f_\tau^*$  denotes the normalized exponential operator for  $\frac{1}{\tau}$ -scaled logits. Below, we derive  $D_{F_\tau^*}(r \parallel s)$  for such  $F_\tau^*$ :

$$\begin{aligned}
D_{F_\tau^*}(s \parallel r) &= F_\tau^*(s) - F_\tau^*(r) - (s - r)^\top \nabla F_\tau^*(r) \\
&= \tau \text{lse}(s/\tau) - \tau \text{lse}(r/\tau) - (s - r)^\top f_\tau^*(r) \\
&= -\tau \left( (s/\tau - \text{lse}(s/\tau)) - (r/\tau - \text{lse}(r/\tau)) \right)^\top f_\tau^*(r) \\
&= \tau f_\tau^*(r)^\top \left( (r/\tau - \text{lse}(r/\tau)) - (s/\tau - \text{lse}(s/\tau)) \right) \\
&= \tau f_\tau^*(r)^\top (\log f_\tau^*(r) - \log f_\tau^*(s)) \\
&= \tau D_{\text{KL}}(f_\tau^*(r) \parallel f_\tau^*(s)) \\
&= \tau D_{\text{KL}}(q \parallel p)
\end{aligned} \tag{12}$$

**Proposition 2.** *The KL divergence between  $p$  and  $q$  in two directions can be expressed as,*

$$\begin{aligned}
D_{\text{KL}}(p \parallel q) &= D_{\text{KL}}(q \parallel p) + \frac{1}{4\tau^2} \text{Var}_{y \sim f^*(a/\tau)}[s(y) - r(y)] - \frac{1}{4\tau^2} \text{Var}_{y \sim f^*(b/\tau)}[s(y) - r(y)] \\
&< D_{\text{KL}}(q \parallel p) + \frac{1}{\tau^2} \|s - r\|_2^2,
\end{aligned} \tag{14}$$

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

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

$$H_{F_\tau^*}(a) = \frac{1}{\tau} (\text{Diag}(f_\tau^*(a)) - f_\tau^*(a) f_\tau^*(a)^\top), \tag{15}$$

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

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

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

$$(s - r)^\top H_{F_\tau^*}(a)(s - r) = \frac{1}{\tau} (s - r)^\top (\text{Diag}(f_\tau^*(a)) - f_\tau^*(a) f_\tau^*(a)^\top)(s - r) \tag{17}$$

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

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

$$\text{and } (s - r)^\top H_{F_\tau^*}(b)(s - r) = \frac{1}{\tau} \text{Var}_{\mathbf{y} \sim f_\tau^*(b)}[s(\mathbf{y}) - r(\mathbf{y})]. \tag{20}$$

Therefore, by combining (19) and (20) with (16) we obtain

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

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

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

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

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

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

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

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).  $\square$

## References

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