

## Mathematical Foundations of Deep Learning (11.80020)

### Assignment 0 (Voluntary exercises) : Problems and Solutions

**Reminder.** For a real random variable  $X$  we call (if existent)

$$\varphi_X(\lambda) := \log \mathbb{E}[e^{\lambda X}] \quad \text{and} \quad \tilde{\varphi}_X(\lambda) := \log \mathbb{E}[e^{\lambda(X - \mathbb{E}X)}]$$

the *logarithmic moment generating function* or shortly *log-moment generating function* and the *centered logarithmic moment generating function* or shortly *centered log-moment generating function*, respectively. Note that  $\varphi_X, \tilde{\varphi}_X: \mathbb{R} \rightarrow \mathbb{R}_{\geq 0} \cup \{+\infty\}$  and we denote their *domains* by

$$\text{dom}(\varphi_X) := \{\lambda \in \mathbb{R} : \varphi_X(\lambda) < \infty\} = \text{dom}(\tilde{\varphi}_X) := \{\lambda \in \mathbb{R} : \tilde{\varphi}_X(\lambda) < \infty\}.$$

We call

$$\varphi_X^*(t) := \sup_{\lambda > 0} \{\lambda t - \varphi_X(\lambda)\} \quad \text{and} \quad \tilde{\varphi}_X^*(t) := \sup_{\lambda > 0} \{\lambda t - \tilde{\varphi}_X(\lambda)\}$$

the *Cramer transform* of  $\varphi_X$  and  $\tilde{\varphi}_X$ , respectively. We call

$$\varphi_X^*(t) := \sup_{\lambda \in \mathbb{R}} \{\lambda t - \varphi_X(\lambda)\} \quad \text{and} \quad \tilde{\varphi}_X^*(t) := \sup_{\lambda \in \mathbb{R}} \{\lambda t - \tilde{\varphi}_X(\lambda)\}$$

the *Legendre transform* (or *Legendre-Fenchel transform* or *convex conjugate*) of  $\varphi_X$  and  $\tilde{\varphi}_X$ , respectively. We call a real random variable *sub-Gaussian* with parameter  $\sigma^2$  (or shortly  $\sigma^2$ -sub-Gaussian) if  $\tilde{\varphi}_X(\lambda) \leq \frac{\sigma^2 \lambda^2}{2}$  for all  $\lambda \in \mathbb{R}$ .

Finally, recall *Hölder's inequality*. For this consider  $p, q \in [1, \infty]$  such that  $\frac{1}{p} + \frac{1}{q} = 1$ , where  $\frac{1}{\infty} := 0$ . For two real random variables  $X$  and  $Y$  it holds that

$$\mathbb{E}[|XY|] \leq \mathbb{E}[|X|^p]^{\frac{1}{p}} \cdot \mathbb{E}[|Y|^q]^{\frac{1}{q}}.$$

**Q1. (Logarithmic moment generating function of a Gaussian)** Compute the centered and non centered log-moment generating function and their Legendre transforms of a Gaussian random variable  $X \sim \mathcal{N}(\mu, \sigma^2)$  with mean  $\mu$  and variance  $\sigma^2$ .

**Solution:** Will be added once the Assignment 1 is handed in.

**Q2. (Stability of sub-Gaussianity)** Let  $X_1$  and  $X_2$  be sub-Gaussian with parameters  $\sigma_1^2$  and  $\sigma_2^2$ , respectively.

(a) Show that  $\alpha X_1$  is sub-Gaussian with parameter  $\alpha^2 \sigma_1^2$  for a constant  $\alpha \in \mathbb{R}$ .

**Solution:** Note that  $\tilde{\varphi}_{\alpha X_1}(\lambda) = \log \mathbb{E}[e^{\lambda(\alpha X_1 - \mathbb{E}[\alpha X_1])}] = \tilde{\varphi}_{X_1}(\alpha \lambda) \leq \frac{\sigma_1^2 \alpha^2 \lambda^2}{2}$ .

(b) If  $X_1$  and  $X_2$  are independent, show that  $X_1 + X_2$  is sub-Gaussian with parameter  $\sigma_1^2 + \sigma_2^2$ .

**Solution:** The independence implies

$$\mathbb{E}[e^{\lambda(X_1 + X_2 - \mathbb{E}[X_1 + X_2])}] = \mathbb{E}[e^{\lambda(X_1 - \mathbb{E}X_1)}] \mathbb{E}[e^{\lambda(X_2 - \mathbb{E}X_2)}].$$

$$\text{Hence, } \tilde{\varphi}_{X_1 + X_2}(\lambda) = \tilde{\varphi}_{X_1}(\lambda) + \tilde{\varphi}_{X_2}(\lambda) \leq \frac{\sigma_1^2 \lambda^2}{2} + \frac{\sigma_2^2 \lambda^2}{2} = \frac{(\sigma_1^2 + \sigma_2^2) \lambda^2}{2}.$$

- (c) Show that in general (without assuming independence), the random variable  $X_1 + X_2$  is sub-Gaussian with parameter  $2(\sigma_1^2 + \sigma_2^2)$ . Next, show that  $X_1 + X_2$  is sub-Gaussian with parameter  $(\sigma_1 + \sigma_2)^2$ , which improves the result.

*Hint:* Use Cauchy-Schwarz and Hölder's inequality respectively.

**Solution:** By Cauchy-Schwarz we have

$$\mathbb{E}[e^{\lambda(X_1+X_2-\mathbb{E}[X_1+X_2])}] \leq \mathbb{E}\left[e^{2\lambda(X_1-\mathbb{E}X_1)}\right]^{\frac{1}{2}} \mathbb{E}\left[e^{2\lambda(X_2-\mathbb{E}X_2)}\right]^{\frac{1}{2}}.$$

Hence, we can estimate

$$\tilde{\varphi}_{X_1+X_2}(\lambda) \leq \frac{\tilde{\varphi}_{X_1}(2\lambda) + \tilde{\varphi}_{X_2}(2\lambda)}{2} \leq \frac{4\sigma_1^2\lambda^2 + 4\sigma_2^2\lambda^2}{4} = \frac{2(\sigma_1^2 + \sigma_2^2)\lambda^2}{2}.$$

By Jensen's inequality with  $p = \frac{\sigma_1+\sigma_2}{\sigma_1}$  and  $q = \frac{\sigma_1+\sigma_2}{\sigma_2}$  we have

$$\mathbb{E}[e^{\lambda(X_1+X_2-\mathbb{E}[X_1+X_2])}] \leq \mathbb{E}\left[e^{\frac{\sigma_1+\sigma_2}{\sigma_1}\lambda(X_1-\mathbb{E}X_1)}\right]^{\frac{\sigma_1}{\sigma_1+\sigma_2}} \cdot \mathbb{E}\left[e^{\frac{\sigma_1+\sigma_2}{\sigma_2}\lambda(X_2-\mathbb{E}X_2)}\right]^{\frac{\sigma_2}{\sigma_1+\sigma_2}}.$$

Hence, we can estimate

$$\begin{aligned} \tilde{\varphi}_{X_1+X_2}(\lambda) &\leq \frac{\sigma_1}{\sigma_1 + \sigma_2} \cdot \tilde{\varphi}_{X_1}\left(\frac{\sigma_1 + \sigma_2}{\sigma_1}\lambda\right) + \frac{\sigma_2}{\sigma_1 + \sigma_2} \cdot \tilde{\varphi}_{X_2}\left(\frac{\sigma_2 + \sigma_2}{\sigma_2}\lambda\right) \\ &\leq \frac{\sigma_1}{\sigma_1 + \sigma_2} \cdot \frac{(\sigma_1 + \sigma_2)^2\lambda^2\sigma_1^2}{2\sigma_1^2} + \frac{\sigma_2}{\sigma_1 + \sigma_2} \cdot \frac{(\sigma_1 + \sigma_2)^2\lambda^2\sigma_2^2}{2\sigma_2^2} \\ &= \frac{(\sigma_1(\sigma_1 + \sigma_2) + \sigma_2(\sigma_1 + \sigma_2)) \cdot \lambda^2}{2} \\ &= \frac{(\sigma_1 + \sigma_2)^2\lambda^2}{2}. \end{aligned}$$

Recall Young's inequality  $2\sigma_1\sigma_2 \leq \sigma_1^2 + \sigma_2^2$  which yields  $(\sigma_1 + \sigma_2)^2 \leq 2(\sigma_1^2 + \sigma_2^2)$ .

**Q3. (Properties of logarithmic moment generating functions)** For this exercise you can assume  $\mathbb{E}X = 0$  in which case  $\varphi_X = \tilde{\varphi}_X$ , but the statements hold in general.

- (a) *Convexity.* Show that  $\varphi_X$  and  $\tilde{\varphi}_X$  are convex functions.

*Hint:* Hölder's inequality

**Solution:** Consider  $\lambda_1, \lambda_2 \in \mathbb{R}$  and  $t \in (0, 1)$ . Then by Hölder's inequality we have

$\mathbb{E}[YZ] \leq \mathbb{E}[Y^{\frac{1}{t}}]^t \mathbb{E}[Z^{\frac{1}{1-t}}]^{1-t}$  for  $Y, Z \geq 0$  and hence we find

$$\begin{aligned} \varphi_X(t\lambda_1 + (1-t)\lambda_2) &= \log\left(\mathbb{E}[e^{(t\lambda_1+(1-t)\lambda_2)X}]\right) = \log\left(\mathbb{E}[e^{t\lambda_1 X} e^{(1-t)\lambda_2 X}]\right) \\ &\leq \log\left(\mathbb{E}[e^{\lambda_1 X}]^t \mathbb{E}[e^{\lambda_2 X}]^{1-t}\right) = t\varphi_X(\lambda_1) + (1-t)\varphi_X(\lambda_2). \end{aligned}$$

- (b) *Semi-continuity.* Show that  $\varphi_X$  and  $\tilde{\varphi}_X$  lower semi-continuous, i.e., if  $\lambda_n \rightarrow \lambda$  for  $n \rightarrow \infty$  then

$$\liminf_{n \rightarrow \infty} \varphi_X(\lambda_n) \geq \varphi_X(\lambda) \quad \text{and} \quad \liminf_{n \rightarrow \infty} \tilde{\varphi}_X(\lambda_n) \geq \tilde{\varphi}_X(\lambda).$$

*Hint:* Fatou's lemma

**Solution:** Using Fatou's lemma we find

$$\liminf_{n \rightarrow \infty} \varphi_X(\lambda_n) = \log\left(\liminf_{n \rightarrow \infty} \mathbb{E}[e^{\lambda_n X}]\right) \geq \log\left(\mathbb{E}[\liminf_{n \rightarrow \infty} e^{\lambda_n X}]\right) = \log \mathbb{E}[e^{\lambda X}] = \varphi_X(\lambda).$$

- (c) *Existence of moments.* Assume that  $\varphi_X(\lambda) < \infty$  or  $\tilde{\varphi}_X(\lambda) < \infty$  for all  $\lambda \in (-\varepsilon, \varepsilon)$  for some  $\varepsilon > 0$ . Show that all moments exist, i.e.,  $\mathbb{E}[|X|^k] < \infty$  for all  $k \in \mathbb{N}$ .

**Solution:** Note that  $e^{\lambda|X|} \leq e^{\lambda X} + e^{-\lambda X}$  and hence for  $\lambda \in (0, \varepsilon)$  we have

$$\mathbb{E}[e^{\lambda|X|}] \leq \varphi_X(\lambda) + \varphi_X(-\lambda) < \infty.$$

In particular, this implies that

$$\frac{\lambda^k \mathbb{E}[|X|^k]}{k!} \leq \mathbb{E}[e^{\lambda|X|}] < \infty.$$

- (d) *Smoothness.* Show that  $\varphi_X$  and  $\tilde{\varphi}_X$  are smooth, i.e., infinitely many times continuously differentiable, functions on the interior of their domains.

*Hint:* It suffices to show that  $e^{\varphi_X}$  and  $e^{\tilde{\varphi}_X}$  are smooth for which you can use the dominated convergence theorem.

**Solution:** It suffices to show smoothness of  $M(\lambda) := e^{\varphi_X(\lambda)} = \mathbb{E}[e^{\lambda X}]$  for which we fix  $\lambda \in \text{int}(\text{dom}(\varphi_X))$ . Recall that integration and differentiation can be exchanged if the is a dominating integrable function (locally) independent of the variable that we consider in the differentiation. In our case this means that we need to find an upper bound  $|\partial_{\lambda'}^{(k)} e^{\lambda' X}| = |X|^k e^{\lambda' X}$  with finite expectation that holds for all  $\lambda' \in (\lambda - \varepsilon, \lambda + \varepsilon)$ . First, note that

$$|X|^k e^{\lambda' X} \leq |X|^k \left( e^{(\lambda+\varepsilon)X} + e^{(\lambda-\varepsilon)X} \right)$$

and hence it remains to show that  $\mathbb{E}[|X|^k e^{(\lambda \pm \varepsilon)X}] < \infty$ . Note that  $(\lambda \pm \varepsilon)p \in \text{dom}(\varphi_X)$  for  $p = (1 - 1/l)^{-1}$  for  $l \in \mathbb{N}$  large enough. Now, we can use Hölder's inequality to show

$$\mathbb{E}[|X|^k e^{(\lambda \pm \varepsilon)X}] \leq \mathbb{E}[|X|^{kl}]^{\frac{1}{l}} \cdot \mathbb{E}[e^{p(\lambda \pm \varepsilon)X}]^{\frac{1}{p}} < \infty$$

since all absolute moments of  $X$  exist.

- (e) *Derivatives.* Show that

$$\varphi'_X(\lambda) = \frac{\mathbb{E}[X e^{\lambda X}]}{M(\lambda)} \quad \text{and} \quad \tilde{\varphi}'_X(\lambda) = \frac{\mathbb{E}[X e^{\lambda(X - \mathbb{E}X)}]}{\widetilde{M}(\lambda)},$$

where  $M(\lambda) := \mathbb{E}[e^{\lambda X}] = e^{\varphi_X(\lambda)}$  and  $\widetilde{M}(\lambda) := \mathbb{E}[e^{\lambda(X - \mathbb{E}X)}] = e^{\tilde{\varphi}_X(\lambda)}$  denote the exponential and centered moment generating functions of  $X$ .

**Solution:** We have seen in the previous exercise that we can exchange integration and differentiation by the dominated convergence theorem. Hence, we can compute

$$\varphi'_X(\lambda) = \partial_{\lambda} \log \mathbb{E}[e^{\lambda X}] = \frac{\mathbb{E}[\partial_{\lambda} e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} = \frac{\mathbb{E}[X e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]}$$

and analogously for  $\tilde{\varphi}_X$ .

- (f) *Cramer transform equals Legendre transform.* Assume that  $\tilde{\varphi}_X$  is finite on  $(-\varepsilon, \varepsilon)$ . Show that

$$\tilde{\varphi}_X^*(t) = \sup_{\lambda > 0} \{\lambda t - \tilde{\varphi}_X(\lambda)\} = \sup_{\lambda \in \mathbb{R}} \{\lambda t - \tilde{\varphi}_X(\lambda)\} = \tilde{\varphi}_X^*(t) \quad \text{for all } t \geq 0.$$

**Solution:** It is clear from the definition that  $\tilde{\varphi}_X^* \leq \tilde{\varphi}_X^*$ . The function  $\tilde{\varphi}_X$  is smooth and in particular continuous and hence

$$\tilde{\varphi}_X^*(t) = \sup_{\lambda > 0} \{\lambda t - \tilde{\varphi}_X(\lambda)\} = \sup_{\lambda \geq 0} \{\lambda t - \tilde{\varphi}_X(\lambda)\}.$$

Note that  $\tilde{\varphi}_X(0) = 0$  and therefore  $\tilde{\varphi}_X^*(0) \geq 0 \cdot t - \tilde{\varphi}_X(0) = 0$ . By Jensen's inequality we have  $\tilde{\varphi}_X(\lambda) \geq \mathbb{E}[\lambda X] = 0$  for all  $\lambda \in \mathbb{R}$  and hence for  $\lambda < 0$  and  $t \geq 0$  we have

$$\lambda t - \tilde{\varphi}_X(\lambda) \leq 0,$$

which shows the desired equality.