

- Solutions should be complete and concisely written. Please, use a separate sheet (or set of sheets) for each problem.
- We will be using Gradescope (<https://www.gradescope.com>) for homework submission (you should have received an invitation) - no paper homework will be accepted. Handwritten solutions are still fine though, just make a good quality scan and upload it to Gradescope.
- You are welcome to discuss problems with your colleagues, but should write and submit your own solution.

## # 1: Properties of exponential families

Recall that an exponential family in canonical form is a class of probability measures on  $\mathbb{R}^n$ , taking the form

$$P_{\theta}(dx) = \frac{1}{Z(\theta)} \exp\{\langle \theta, T(x) \rangle\} \nu(dx), \quad (1)$$

where  $\nu(dx)$  is a reference measure on  $\mathbb{R}^n$ . For the purpose of this problem, you can assume that  $P_{\theta}$  has a density with respect to the Lebesgue measure on  $\mathbb{R}^n$ , which therefore can be written as

$$p_{\theta}(x) = \frac{1}{Z(\theta)} \exp\{\langle \theta, T(x) \rangle\} h(x), \quad (2)$$

where  $h : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$  is a measurable function. Alternatively, you can assume that  $P_{\theta}$  is supported on  $\mathbb{Z}^d$ , with probability mass function of the form (2). Recall that the log partition function is defined as  $\phi(\theta) = \log Z(\theta)$ , which is finite for  $\theta \in \Theta_N$  (the natural parameter space). (In the following, you are not required to justify the exchange of order of derivative and integrals.)

- (a) Prove that  $\Theta_N$  is convex and  $\phi : \Theta_N \rightarrow \mathbb{R}$  is a convex function.
- (b) Prove the following identities hold for  $\theta \in \Theta_N^{\circ}$  (the interior of  $\Theta_N$ )

$$\frac{\partial \phi}{\partial \theta_i}(\theta) = E_{\theta}\{T_i(\mathbf{X})\}, \quad (3)$$

$$\frac{\partial^2 \phi}{\partial \theta_i \partial \theta_j}(\theta) = \text{Cov}_{\theta}\{T_i(\mathbf{X}); T_j(\mathbf{X})\}. \quad (4)$$

- (c) Assume that  $\mathbf{x} \mapsto h(\mathbf{x})$ ,  $\mathbf{x} \mapsto T(\mathbf{x})$  are differentiable. Prove that the following identity (Stein's identity) hold for any differentiable function  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  such that both sides make sense:

$$E_{\theta}\left\{\left[\frac{1}{h(\mathbf{x})} \frac{\partial h}{\partial x_i}(\mathbf{x}) + \langle \theta, \frac{\partial T}{\partial x_i}(\mathbf{x}) \rangle\right] g(\mathbf{x})\right\} + E_{\theta}\left\{\frac{\partial g}{\partial x_i}(\mathbf{x})\right\} = 0 \quad (5)$$

- (d) Assume that  $p$  is a multivariate Gaussian density, namely

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} \det(\Sigma)^{1/2}} \exp\left\{-\frac{1}{2} \langle (\mathbf{x} - \mu) \Sigma^{-1} (\mathbf{x} - \mu) \rangle\right\}. \quad (6)$$

Show that Stein's identity in this case reduces to

$$E\{(\mathbf{x} - \mu) g(\mathbf{x})\} = \Sigma E\{\nabla g(\mathbf{x})\}. \quad (7)$$

## # 2: Exercises on sufficient statistics

- (a) Consider a statistical model composed of  $k$  probability distributions. Namely

$$\mathcal{P} = \{p_1, p_2, \dots, p_k\}, \quad (8)$$

where  $p_\ell$  are densities on  $\mathbb{R}^n$  (we identify the probability distribution with its density).

Show that there exists a set of  $k - 1$  sufficient statistics.

- (b) Let  $\boldsymbol{\theta} = (\theta_1, \theta_2) \in \mathbb{R}^2$  with  $\theta_1 < \theta_2$  and define  $\mathsf{P}_{\boldsymbol{\theta}}$  to be the distribution of  $n$  i.i.d. random variables  $X_1, \dots, X_n \sim \text{Unif}([\theta_1, \theta_2])$ . Let  $x_{\min} = \min(x_1, \dots, x_n)$ , and  $x_{\max} = \max(x_1, \dots, x_n)$ . Prove that  $(x_{\min}, x_{\max})$  is a sufficient statistics for the model  $\mathcal{P} = (\mathsf{P}_{\boldsymbol{\theta}})$ .
- (c) Consider the Gaussian linear model. Namely, for a fixed design matrix  $\mathbf{A} \in \mathbb{R}^{n \times d}$ , we have  $\mathsf{P}_{\boldsymbol{\theta}} = \mathsf{N}(\mathbf{A}\boldsymbol{\theta}, \sigma^2 \mathbf{I}_n)$ ,  $\boldsymbol{\theta} \in \Theta = \mathbb{R}^d$ . Show that there exists a sufficient statistic of dimension  $d$ .

## # 3: Optimal linear estimation in heteroscedastic Gaussian model

Assume  $\sigma_1, \dots, \sigma_d > 0$  to be known, and consider the statistical model  $\mathsf{P}_{\boldsymbol{\theta}} = \mathsf{N}(\boldsymbol{\theta}\mathbf{1}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ , and  $\boldsymbol{\theta} \in \Theta = \mathbb{R}$  (with  $\mathbf{1}$  denoting the all-ones vector). In other words,  $X_i = \theta + \sigma_i G_i$  where  $(G_i)_{i \leq d} \sim_{\text{iid}} \mathsf{N}(0, 1)$ . (Here  $\langle \mathbf{u}, \mathbf{v} \rangle = \sum_{i=1}^m u_i v_i$  denotes the usual scalar product of  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^m$ .)

- (a) Show that there exists a sufficient statistics of the form  $\ell(\mathbf{x}) = \langle \mathbf{c}, \mathbf{x} \rangle$ , where  $\mathbf{c} \in \mathbb{R}^d$ , and determine the vector  $\mathbf{c}$ .
- (b) Using the result at the previous point, determine the optimal linear estimator  $\hat{\theta}(\mathbf{x}) = \langle \mathbf{a}, \mathbf{x} \rangle$ , with respect to the square loss  $L(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$ . Here optimality is to be understood in minimax sense, that is we want to minimize  $R_{\mathsf{M}}(\hat{\theta})$  among all linear estimators.
- (c) Generalize the above to the case of general correlated Gaussian noise (i.e.  $\boldsymbol{\Sigma}$  not necessarily diagonal, but strictly positive definite).