# Deep Learning Quiz

Date: 9/22/2022. Duration: 50 minutes

1. Answer True or False in the following statements about the Gaussian distribution $\mathcal{N}$:

   (a) If $x \sim \mathcal{N}$, then $ax + b \sim \mathcal{N}$ for any constants $a, b \in \mathbb{R}$ (5%)

   (b) If $z = x + y$, where $x, y \sim \mathcal{N}$, then $z \sim \mathcal{N}$. (5%)

2. Given $N$ i.i.d samples $X^{\in N \times D} = [x^{(1)}, \dots, x^{(N)}]^T$ of a random variable $\mathbf{x}$, the Principal Components Analysis (PCA) finds $K$ orthonormal vector $W = [w^{(1)}, \dots, w^{(K)}]$ such that the transformed variable $\mathbf{z} = W^T \mathbf{x}$ has the most "spread out" attributes, i.e., each attribute $z_i = w^{(i)T}\mathbf{x}$ has the maximum variance $\text{Var}(z_i)$. Now consider the problem of finding $w^{(1)}$:

   (a) Assuming that $\mathbf{x}$ has zero mean, show that $\sigma_{z1}^2 = \frac{1}{N} w^{(1)T} X^T X w^{(1)}$. (10%)

   (b) Use the Rayleigh's Quotient to explain that the optimal $w^{(1)}$ is given by the eigenvector of $X^T X$ corresponding to the largest eigenvalue. (10%)

3. Consider a situation where a doctor wants to inference if a patient is having either the disease $y^{(1)}$ or $y^{(2)}$ by examining the patient's symptoms $\mathbf{x}$. Explain why the Bayes' rule,

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)},$$

   can make the inference easier. (10%)

4. Give an example of two distributions $P$ and $Q$ to show that the Kullback-Leibler (KL) Divergence

   $D_{KL}(P||Q) = E_{x \sim P}[\log \frac{P(x)}{Q(x)}]$ is asymmetric, i.e., $D_{KL}(P||Q) \neq D_{KL}(Q||P)$. (10%)

5. Consider a continuous, differentiable function $f: \mathbb{R}^d \to \mathbb{R}$ and an input point $\mathbf{a} \in \mathbb{R}^d$.

   (a) For any direction $\mathbf{u}$ in the input space, show that the directional derivative of $f$ at $\mathbf{a}$ along $\mathbf{u}$ equals to $\nabla f(\mathbf{a})^T \mathbf{u}$. (10% Hint: the directional derivative of $f$ at $\mathbf{a}$ along $\mathbf{u}$ is the derivative of function $f(\mathbf{a} + \varepsilon\mathbf{u})$ with respect to $\varepsilon$, evaluated at $\varepsilon = 0$.)

   (b) What is the direction in the input space that leads to the steepest decent of $f$ starting from $\mathbf{a}$, i.e., what is the solution of $\text{argmin}_{\mathbf{u}, \|\mathbf{u}\|=1} \nabla f(\mathbf{a})^T \mathbf{u}$? (10%)

6. Given a quadratic function $f: \mathbb{R}^2 \to \mathbb{R}$, $f(x) = \frac{1}{2}x^T A x - b^T x + c$, where $A \in \mathbb{R}^{2 \times 2}$ is symmetric. Explain why the problem

$$\text{argmin}_x f(x)$$

   is hard to solve by Gradient Descent algorithm when $A$ is ill-conditioned (i.e., when the condition number $\varkappa(A) = \max_{i,j} |\frac{\lambda_i}{\lambda_j}|$ is large). (10%)

7. Given a vector $x$, let $z = x - \max_i x_i 1$. When you implement the line $c = \log(softmax(z)_i)$ for some $i > 0$ in a computer program which stores $z_i$ as a float, what numerical issues may occur? (10%) How to walk around these issues in your implementations? (10%)

8. Consider a constrained optimization problem:

$$\min_x f(x)$$

$$\text{subject to } x \in \{x: g^{(i)}(x) \leq 0, h^{(j)}(x) = 0\}_{i,j},$$

for some positive integers $i$ and $j$. Explain why the following unconstrained problem:

$$\min_x \max_{\alpha,\beta,\alpha \geq 0} f(x) + \sum_i \alpha_i g^{(i)}(x) + \sum_j \beta_j h^{(j)}(x)$$

gives the same optimal solution. (10%)