**Homework 3 (Calculus) Problem 1**

Honor statement:

Read:

**Solution:** Complete        ∎

Read:

**Solution:** Complete                                                 ■

Read:

**Solution:** Complete              ∎

Let $x$ be the $K \times 1$ vector output of the last layer of a xNN and $e = \text{crossEntropy}(p^*, \text{softMax}(x))$ be the error where $p^*$ is a $K \times 1$ vector with a 1 in position $k^*$ representing the correct class and 0s elsewhere. Derive $\partial e / \partial x$

**Solution:**

$$\frac{\partial e}{\partial x} = ?  \tag{1}$$

■

**Proof:** First, note that error function $e$ is given by the cross entropy of a softmax, the typical error function chosen for **classification** networks. This establishes a $f(g(x))$ relationship wherein we must use the chain rule to differentiate $e$. Specifically, we will need to compute the following to apply the chain rule

$$\frac{d}{dx}\text{softMax}(x)  \qquad\qquad  \frac{\partial}{\partial p}\text{crossEntropy}(p^*, p)  \tag{2}$$

We will start with cross entropy. Recall that cross entropy is given by

$$\text{crossEntropy}(p^*, p) = -\sum_{x_i \in x} p^*(x_i) \log p(x_i)  \tag{3}$$

Differentiating cross entropy for each of $x_i \in x$ produces a gradient vector. Recognizing that differentiating the summation in cross entropy with respect to a single $x_i$ will produce a single nonzero term, we can say

$$\nabla_x\left(-\sum_{x_i \in x} p^*(x_i)\log p(x_i)\right) = \begin{pmatrix} -\frac{\partial}{\partial x_1}p^*(x_1)\log p(x_1) \\ -\frac{\partial}{\partial x_2}p^*(x_2)\log p(x_2) \\ \vdots \\ -\frac{\partial}{\partial x_k}p^*(x_k)\log p(x_k) \\ \vdots \\ -\frac{\partial}{\partial x_N}p^*(x_N)\log p(x_N) \end{pmatrix}  \tag{4}$$

$$= \begin{pmatrix} 0 \\ 0 \\ \vdots \\ -\frac{1}{p(x_k)} \\ \vdots \\ 0 \end{pmatrix} \leftarrow k  \tag{5}$$

So our cross entropy gradient for a one hot vector $\boldsymbol{p}^*$ with nonzero position $k$ is 0 everywhere except at position $k$.

Now we compute our next required derivative: softmax The softmax function is given by

$$\text{softmax}(\boldsymbol{x}) = \frac{1}{\sum_{k=1}^{K} e^{x_k}} \begin{pmatrix} e^{x_0} \\ e^{x_1} \\ \vdots \\ e^{x_K} \end{pmatrix} \tag{6}$$

We need to derive the Jacobian of this function in denominator notation. $\qquad \square$

(a) What is the arithmetic intensity for matrix matrix multiplication with sizes $M, N, K$.