

Question 1

(a)

Since $\sigma_i(\mathbf{x}) > 0$ for all x, i , and its value is dependent on all other $j \neq i$, every entry of the Jacobian matrix is non-zero.

In order to get rid of the exponential functions, we can try to use the logarithm of the softmax function:

$$\begin{aligned}\frac{\partial \log \sigma_i(\mathbf{x})}{\partial x_j} &= \frac{1}{\sigma_i(\mathbf{x})} \frac{\partial \sigma_i(\mathbf{x})}{\partial x_j} \\ \frac{\partial \sigma_i(\mathbf{x})}{\partial x_j} &= \sigma_i(\mathbf{x}) \frac{\partial \log \sigma_i(\mathbf{x})}{\partial x_j}\end{aligned}$$

$$\begin{aligned}\log \sigma_i(\mathbf{x}) &= \log \left(\frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} \right) \\ &= x_i - \log \left(\sum_{j=1}^n \exp(x_j) \right)\end{aligned}$$

In the following we use:

$$\frac{\partial x_i}{\partial x_z} = \begin{cases} 1 & \text{if } i = z \\ 0 & \text{if } i \neq z \end{cases}$$

$$\begin{aligned}\frac{\partial \log \sigma_i(\mathbf{x})}{\partial x_j} &= \frac{\partial x_i}{\partial x_j} - \frac{\partial \log \left(\sum_{j=1}^n \exp(x_j) \right)}{\partial x_j} \\ &= \mathbb{1}_{i=j} - \frac{\partial \log \left(\sum_{j=1}^n \exp(x_j) \right)}{\partial x_j} \\ &= \mathbb{1}_{i=j} - \frac{1}{\sum_{j=1}^n \exp(x_j)} \left(\frac{\partial}{\partial x_j} \sum_{j=1}^n \exp(x_j) \right) \\ &= \mathbb{1}_{i=j} - \frac{\exp(x_j)}{\sum_{j=1}^n \exp(x_j)} \\ &= \mathbb{1}_{i=j} - \sigma_j(\mathbf{x})\end{aligned}$$

Finally, convert back to the original derivative:

$$\begin{aligned}\frac{\partial \sigma_i(\mathbf{x})}{\partial x_j} &= \sigma_i(\mathbf{x}) \frac{\partial \log \sigma_i(\mathbf{x})}{\partial x_j} \\ &= \sigma_i(\mathbf{x}) \cdot (\mathbb{1}_{i=j} - \sigma_j(\mathbf{x}))\end{aligned}$$

$$\begin{aligned}
D_{\mathbf{x}}\sigma(\mathbf{x}) &= \begin{pmatrix} \frac{\partial\sigma_1(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial\sigma_1(\mathbf{x})}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial\sigma_n(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial\sigma_n(\mathbf{x})}{\partial x_n} \end{pmatrix} \\
&= \begin{pmatrix} \sigma_1(\mathbf{x}) \cdot (1 - \sigma_1(\mathbf{x})) & \cdots & -\sigma_1(\mathbf{x}) \cdot \sigma_n(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ -\sigma_n(\mathbf{x}) \cdot \sigma_1(\mathbf{x}) & \cdots & \sigma_n(\mathbf{x}) \cdot (1 - \sigma_n(\mathbf{x})) \end{pmatrix}
\end{aligned}$$

So especially, the diagonal entries are:

$$\frac{\partial\sigma_i(\mathbf{x})}{\partial x_i} = \sigma_i(\mathbf{x}) \cdot (1 - \sigma_i(\mathbf{x}))$$

And the off-diagonal entries are:

$$\frac{\partial\sigma_i}{\partial x_j} = -\sigma_i(\mathbf{x}) \cdot \sigma_j(\mathbf{x})$$

And the matrix is symmetric. Thus:

$$\begin{aligned}
\frac{\partial\sigma_i}{\partial x_j} &= -\sigma_i(\mathbf{x}) \cdot \sigma_j(\mathbf{x}) \\
&= -\sigma_j(\mathbf{x}) \cdot \sigma_i(\mathbf{x}) \\
&= \frac{\partial\sigma_j}{\partial x_i}
\end{aligned}$$

(b)

$$\begin{aligned}
\mathbf{z} &= \mathbf{v} \cdot D_{\mathbf{x}}\sigma(\mathbf{x}) \\
&= (v_1 \dots v_n) \cdot \begin{pmatrix} \sigma_1(\mathbf{x}) \cdot (1 - \sigma_1(\mathbf{x})) & \cdots & -\sigma_1(\mathbf{x}) \cdot \sigma_n(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ -\sigma_n(\mathbf{x}) \cdot \sigma_1(\mathbf{x}) & \cdots & \sigma_n(\mathbf{x}) \cdot (1 - \sigma_n(\mathbf{x})) \end{pmatrix} \\
&= \begin{pmatrix} v_1 \cdot \sigma_1(\mathbf{x}) \cdot (1 - \sigma_1(\mathbf{x})) + \cdots + v_n \cdot -\sigma_n(\mathbf{x}) \cdot \sigma_1(\mathbf{x}) \\ \vdots \\ v_1 \cdot -\sigma_1(\mathbf{x}) \cdot \sigma_n(\mathbf{x}) + \cdots + v_n \cdot \sigma_n(\mathbf{x}) \cdot (1 - \sigma_n(\mathbf{x})) \end{pmatrix}^{\top} \\
&= \begin{pmatrix} \sigma_1(\mathbf{x}) \cdot (v_1 \cdot (1 - \sigma_1(\mathbf{x})) - v_2 \cdot \sigma_2(\mathbf{x}) - \cdots - v_n \cdot \sigma_n(\mathbf{x})) \\ \vdots \\ \sigma_n(\mathbf{x}) \cdot (v_1 \cdot \sigma_1(\mathbf{x}) - v_{n-1} \cdot \sigma_{n-1}(\mathbf{x}) - \cdots + v_n \cdot (1 - \sigma_n(\mathbf{x}))) \end{pmatrix}^{\top} \\
&= \begin{pmatrix} \sigma_1(\mathbf{x}) \cdot (v_1 - v_1 \cdot \sigma_1(\mathbf{x}) - v_2 \cdot \sigma_2(\mathbf{x}) - \cdots - v_n \cdot \sigma_n(\mathbf{x})) \\ \vdots \\ \sigma_n(\mathbf{x}) \cdot (v_1 \cdot \sigma_1(\mathbf{x}) - v_{n-1} \cdot \sigma_{n-1}(\mathbf{x}) - \cdots + v_n - v_n \cdot \sigma_n(\mathbf{x})) \end{pmatrix}^{\top} \\
&= \begin{pmatrix} \sigma_1(\mathbf{x}) \cdot (v_1 - \mathbf{v} \cdot \sigma(\mathbf{x})^{\top}) \\ \vdots \\ \sigma_n(\mathbf{x}) \cdot (v_n - \mathbf{v} \cdot \sigma(\mathbf{x})^{\top}) \end{pmatrix}^{\top}
\end{aligned}$$

(c)

$$\begin{aligned}
\frac{\partial l(\mathbf{z}, \mathbf{t})}{z_j} &= -\frac{\partial}{\partial z_j} \sum_{i=1}^n t_i \cdot \log(z_i) \\
&= -\sum_{i=1}^n t_i \cdot \frac{\partial}{\partial z_j} \log(z_i) \\
&= -\sum_{i=1}^n \frac{t_i}{z_i} \cdot \frac{\partial z_i}{\partial z_j} \\
&= -\frac{t_i}{z_i} \cdot \frac{\partial}{\partial z_j} \sum_{i=1}^n z_i \\
&= -\frac{t_i}{z_i}
\end{aligned}$$

$$D_{\mathbf{z}} l(\mathbf{z}, \mathbf{t}) = \begin{pmatrix} -\frac{t_1}{z_1} \\ \vdots \\ -\frac{t_n}{z_n} \end{pmatrix}^\top$$

(d)

If one of the terms $z_i = 0$, then $D_{\mathbf{z}} l(\mathbf{z}, \mathbf{t})$ is not computable.

We have already observed that $\sigma_i(\mathbf{x}) > 0$. Thus, $z_i = 0$ can only occur when the following is satisfied:

$$v_j = \mathbf{v} \cdot \sigma(\mathbf{x})^\top$$