

2. (b)

$$\begin{aligned} J &= - \sum_{i=1}^V y_i \ln \hat{y}_i = - \{\mathbf{y}\}^\top \ln(\text{softmax}(\{\boldsymbol{\theta}\})) \\ &= - \langle \mathbf{y} \rangle \left(\{\boldsymbol{\theta}\} - \ln \sum_{i=1}^V e^{\theta_i} \right). \end{aligned}$$

$$\begin{aligned} \frac{\partial J}{\partial \langle \boldsymbol{\theta} \rangle} &= - \langle \mathbf{y} \rangle \left([\mathbf{I}] - \frac{1}{\sum e^{\theta_i}} \cdot \sum_{i=1}^V \frac{\partial e^{\theta_i}}{\partial \langle \boldsymbol{\theta} \rangle} \right) \\ &= - \langle \mathbf{y} \rangle \left([\mathbf{I}] - \frac{e^{\langle \boldsymbol{\theta} \rangle}}{\sum e^{\theta_i}} \right) \\ &= - \langle \mathbf{y} \rangle ([\mathbf{I}] - \langle \hat{\mathbf{y}} \rangle) \\ &= \langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle. \end{aligned}$$

3. (a)

$$\begin{aligned} \frac{\partial J}{\partial \langle \mathbf{v}_c \rangle} &= \frac{\partial J}{\partial \langle \boldsymbol{\theta} \rangle} \frac{\partial \{\boldsymbol{\theta}\}}{\partial \langle \mathbf{v}_c \rangle} \\ &= (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle) \frac{\partial \{\boldsymbol{\theta}\}}{\partial \langle \mathbf{v}_c \rangle}. \end{aligned}$$

Since

$$\{\boldsymbol{\theta}\} = \begin{Bmatrix} \langle \mathbf{u}_1 \rangle \{\mathbf{v}_c\} \\ \langle \mathbf{u}_2 \rangle \{\mathbf{v}_c\} \\ \vdots \\ \langle \mathbf{u}_V \rangle \{\mathbf{v}_c\} \end{Bmatrix} = [\{\mathbf{u}_1\} \quad \{\mathbf{u}_2\} \quad \cdots \quad \{\mathbf{u}_V\}]^\top \{\mathbf{v}_c\},$$

therefore,

$$\begin{aligned} \frac{\partial J}{\partial \langle \mathbf{v}_c \rangle} &= (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle) [\mathbf{U}]^\top. \\ \frac{\partial J}{\partial \{\mathbf{v}_c\}} &= [\mathbf{U}] (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle). \end{aligned}$$

3. (b)

$$\begin{aligned}
\frac{\partial J}{\partial [\mathbf{U}]} &= (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle) \frac{\partial \{\boldsymbol{\theta}\}}{\partial [\mathbf{U}]} \\
&= \begin{bmatrix} (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle) \begin{bmatrix} \frac{\partial \{\mathbf{u}_1\}^\top \{\mathbf{v}_c\}}{\partial \langle \mathbf{u}_1 \rangle} \\ \vdots \\ \frac{\partial \{\mathbf{u}_1\}^\top \{\mathbf{v}_c\}}{\partial \langle \mathbf{u}_V \rangle} \end{bmatrix} \\ (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle) \begin{bmatrix} \frac{\partial \{\mathbf{u}_2\}^\top \{\mathbf{v}_c\}}{\partial \langle \mathbf{u}_1 \rangle} \\ \vdots \\ \frac{\partial \{\mathbf{u}_2\}^\top \{\mathbf{v}_c\}}{\partial \langle \mathbf{u}_V \rangle} \end{bmatrix} \\ \vdots \\ (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle) \begin{bmatrix} \frac{\partial \{\mathbf{u}_V\}^\top \{\mathbf{v}_c\}}{\partial \langle \mathbf{u}_1 \rangle} \\ \vdots \\ \frac{\partial \{\mathbf{u}_V\}^\top \{\mathbf{v}_c\}}{\partial \langle \mathbf{u}_V \rangle} \end{bmatrix} \end{bmatrix} \\
&= \begin{bmatrix} (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle) \begin{bmatrix} \langle \mathbf{v}_c \rangle \\ \langle \mathbf{0} \rangle \\ \vdots \\ \langle \mathbf{0} \rangle \\ \langle \mathbf{0} \rangle \\ \langle \mathbf{v}_c \rangle \\ \vdots \\ \langle \mathbf{0} \rangle \end{bmatrix} \\ \vdots \\ (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle) \begin{bmatrix} \langle \mathbf{0} \rangle \\ \vdots \\ \langle \mathbf{v}_c \rangle \end{bmatrix} \end{bmatrix} \\
&= \begin{bmatrix} (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle)_{*1} \cdot \langle \mathbf{v}_c \rangle \\ (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle)_{*2} \cdot \langle \mathbf{v}_c \rangle \\ \vdots \\ (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle)_{*V} \cdot \langle \mathbf{v}_c \rangle \end{bmatrix} \\
&= \{\mathbf{v}_c\} (\langle \hat{\mathbf{y}} \rangle - \langle \mathbf{y} \rangle).
\end{aligned}$$