## A Gradient Calculation

The gradients will be derived based on the model with scaled dot-product attention as discussed in the main paper. We consider the case of a single training instance for simplicity and clarity. The gradient on the entire training set can be calculated as the average of the gradients on each training instance.

### A.1 Gradients of $s$, $h$, $W$

$$
\begin{aligned}
\frac{\partial \ell}{\partial s} &= -y\beta, \\
\frac{\partial \ell}{\partial W} &= \frac{\partial \ell}{\partial s}\frac{\partial s}{\partial W} \\
&= -y\beta h, \\
\frac{\partial \ell}{\partial h} &= \frac{\partial \ell}{\partial s}\frac{\partial s}{\partial h} \\
&= -y\beta W,
\end{aligned}
\tag{1}
$$

where $\beta = 1 - \sigma(ys) = \sigma(-ys)$.

### A.2 Gradient of $e$

In our model with scaled dot-product attention, the input of the attention layer consists of word embeddings. The input can be generalized to the output of an affine layer. Here, we will use $h_i$ instead of $e_i$ to denote the general representation of the token at the $i$-th position in the instance. The partial derivative of $h$ with respect to the representation of the $i$-th token $h_i$ can be calculated as:

$$
\begin{aligned}
\frac{\partial h}{\partial h_i} &= \frac{\partial \sum_j \alpha_j h_j}{\partial h_i} \\
&= \sum_j \frac{\partial \alpha_j}{\partial h_i} h_j^\top + \alpha_i \mathbf{I},
\end{aligned}
\tag{2}
$$

where $\mathbf{I} \in \mathbb{R}^{d \times d}$ is an identity matrix. Note that, $h_i$ can be the embedding of the $i$-th token, or the $i$-th output of an affine layer. The partial derivative of $\alpha_j$ with respect to $h_i$ can be represented as:

$$
\frac{\partial \alpha_j}{\partial h_i} = \frac{\partial \alpha_j}{\partial a_i}\frac{\partial a_i}{\partial h_i} = \alpha_j(\delta_{ji} - \alpha_i)\frac{\partial a_i}{\partial h_i}, \tag{3}
$$

where $\delta_{ji}$ is the Kronecker delta function:

$$
\delta_{ji} = \begin{cases} 1 & j = i \\ 0 & else. \end{cases}
\tag{4}
$$

Referring to Equation 3 and 4 , we can obtain the partial derivative of $h$ with respect to $h_i$ as follows:

$$
\begin{aligned}
\frac{\partial h}{\partial h_i} &= \sum_j \alpha_j(\delta_{ji} - \alpha_i)\frac{\partial a_i}{\partial h_i} h_j^\top + \alpha_i \mathbf{I} \\
&= \alpha_i \left[ \frac{\partial a_i}{\partial h_i}(h_i - h)^\top + \mathbf{I} \right].
\end{aligned}
\tag{5}
$$

On the model with dot-product attention, the attention score and corresponding partial derivative will be written as:

$$
\begin{aligned}
a_i &= \frac{h_i^\top V}{\lambda}, \\
\frac{\partial a_i}{\partial h_i} &= \frac{V}{\lambda},
\end{aligned}
\tag{6}
$$

where $V \in \mathbb{R}^d$ is the context vector. With Equation 5 and 6, we can have:

$$
\begin{aligned}
\frac{\partial h}{\partial h_i} &= \sum_j \alpha_j(\delta_{ji} - \alpha_i)\frac{\partial a_i}{\partial h_i} h_j^\top + \alpha_i \mathbf{I} \\
&= \alpha_i \left[ \frac{V(h_i - h)^\top}{\lambda} + \mathbf{I} \right].
\end{aligned}
\tag{7}
$$

With Equation 1 and 7, we can obtain the partial derivative of the loss with respect to $h_i$:

$$
\begin{aligned}
\frac{\partial \ell}{\partial h_i} &= \frac{\partial h}{\partial h_i}\frac{\partial \ell}{\partial h} \\
&= -y\beta\alpha_i \left[ \frac{V(h_i - h)^\top}{\lambda} + \mathbf{I} \right] W.
\end{aligned}
\tag{8}
$$

As the inputs of the attention layer are embeddings directly, given a token $e$, namely $h_i = e$, the partial derivative on the entire training set can be represented as:

$$
\frac{\partial \ell}{\partial e} = -\frac{1}{m}\sum_{(t,j):e_j^{(t)}\equiv e} y^{(t)}\beta^{(t)}\alpha_j^{(t)}\left[ \frac{V(e - h^{(t)})^\top}{\lambda} + \mathbf{I} \right]W,
\tag{9}
$$

where $(t, j) : e_j^{(t)} \equiv e$ means we are selecting such tokens from the $t$-th instance at the $j$-th position that are exactly $e$, and $\alpha_j^{(t)}$ is the attention weight for that $j$-th token in the selected $t$-th instance.

### A.3 Gradient of $V$

The partial derivative of $\ell$ with respect to $V$ can be written as:

$$
\begin{aligned}
\frac{\partial \ell}{\partial V} &= \frac{\partial h}{\partial V}\frac{\partial \ell}{\partial h} \\
&= \sum_j \frac{\partial \alpha_j h_j}{\partial V}\frac{\partial \ell}{\partial h} \\
&= \sum_j \frac{\partial \alpha_j}{\partial V}h_j^\top\frac{\partial \ell}{\partial h},
\end{aligned}
\tag{10}
$$

note that the attention weight $\alpha_j$ is the function of $V$. The partial derivative of $\alpha_j$ with respect to $V$ is:

$$\frac{\partial \alpha_j}{\partial V} = \sum_i \frac{\partial \alpha_j}{\partial a_i}\frac{\partial a_i}{\partial V}$$

$$= \sum_i \alpha_j(\delta_{ji} - \alpha_i)\frac{h_i}{\lambda} \qquad (11)$$

$$= \frac{\alpha_j}{\lambda}(h_j - h).$$

Substitute Equation 11 into Equation 10, and we will obtain the following:

$$\frac{\partial \ell}{\partial V} = -\frac{y\beta}{\lambda}\sum_j \alpha_j(h_j - h)h_j^\top W$$

$$= -\frac{y\beta}{\lambda}[\sum_j \alpha_j h_j h_j^\top - hh^\top]W \qquad (12)$$

$$= -\frac{y\beta}{\lambda}\sum_j \alpha_j h_j(h_j - h)^\top W.$$

For the simple model in our paper, $h_j = e_j$, the gradient of $V$ on the entire training set can be calculated as:

$$\frac{\partial \ell}{\partial V}=$$

$$-\frac{1}{m\lambda}\sum_{t=1}^m y^{(t)}\beta^{(t)}\sum_j \alpha_j^{(t)}e_j^{(t)}\left(e_j^{(t)} - h^{(t)}\right)^\top W$$

$$= -\frac{1}{m\lambda}\sum_{t=1}^m y^{(t)}\beta^{(t)}\sum_j \alpha_j^{(t)}e_j^{(t)}\left(s_j^{(t)} - s^{(t)}\right), \qquad (13)$$

where $(t, j) : e_j^{(t)} \equiv e$ means we are selecting such tokens from the $t$-th instance at the $j$-th position that are exactly $e$, and $\alpha_j^{(t)}$ is the attention weight for that $j$-th token in the selected $t$-th instance.

## B  Model with an Affine Input Layer

Let us consider a model with an affine input layer. The affine layer is added between the embedding layer and the attention layer. The variables will be described as:

$$h_i = W^e e_i, \qquad (14)$$

$$a_i = \frac{h_i^\top V}{\lambda}, \qquad (15)$$

$$s_i = h_i^\top W, \qquad (16)$$

where $W^e \in R^{d\times d}$ is the weight matrix of the affine layer. The gradients of the parameters $e$, $W^e$

and $V$ on the entire dataset will be calculated as:

$$\frac{\partial \ell}{\partial e} = -\frac{1}{m}\sum_{(t,j):e_j^{(t)}\equiv e} y^{(t)}\beta^{(t)}\alpha_j^{(t)}(s_e - s^{(t)})\frac{(W^e)^\top V}{\lambda}$$

$$- \frac{1}{m}\sum_{(t,j):e_j^{(t)}\equiv e} y^{(t)}\beta^{(t)}\alpha_j^{(t)}(W^e)^\top W,$$

$$\frac{\partial \ell}{\partial W^e} = -\frac{1}{m}\sum_{t=1}^m\sum_j y^{(t)}\beta^{(t)}\alpha_j^{(t)}(s_j^{(t)} - s^{(t)})\frac{V(e_j^{(t)})^\top}{\lambda}$$

$$- \frac{1}{m}\sum_{t=1}^m\sum_j y^{(t)}\beta^{(t)}\alpha_j^{(t)} W(e_j^{(t)})^\top,$$

$$\frac{\partial \ell}{\partial V} = -\frac{1}{m}\sum_{t=1}^m\sum_j y^{(t)}\beta^{(t)}\alpha_j^{(t)}(s_j^{(t)} - s^{(t)})\frac{h_j^{(t)}}{\lambda}. \qquad (17)$$

The update of the polarity score can be represented as:

$$\frac{ds_e}{d\tau} = (\frac{dh_e}{d\tau})^\top W + h_e^\top \frac{dW}{d\tau}$$

$$= (\frac{dW^e e}{d\tau})^\top W + h_e^\top \frac{dW}{d\tau}$$

$$= (\frac{de}{d\tau})^\top (W^e)^\top W + e^\top (\frac{dW^e}{d\tau})^\top W + h_e^\top \frac{dW}{d\tau}. \qquad (18)$$

The update of $s_e$ can be described as below:

$$\frac{ds_e}{d\tau} = \frac{1}{m}\sum_{(t,j):e_j^{(t)}\equiv e} y^{(t)}\beta^{(t)}\alpha_j^{(t)}||W^\top W^e||_2^2$$

$$+ \frac{1}{m}\sum_{t=1}^m y^{(t)}\beta^{(t)} h_e^\top h^{(t)}$$

$$+ \frac{1}{m}\sum_{t=1}^m\sum_j y^{(t)}\beta^{(t)}\alpha_j^{(t)} e^\top e_j^{(t)}||W||_2^2$$

$$+ \frac{1}{m\lambda}\sum_{(t,j):e_j^{(t)}\equiv e} y^{(t)}\beta^{(t)}\alpha_j^{(t)}(s_e - s^{(t)})W'^\top V'$$

$$+ \frac{1}{m\lambda}\sum_{t=1}^m\sum_j y^{(t)}\beta^{(t)}\alpha_j^{(t)}(s_j^{(t)} - s^{(t)})e^\top e_j^{(t)} Q. \qquad (19)$$

where $W' = (W^e)^\top W$, $V' = (W^e)^\top V$, $Q = V^T W$.

The update of the attention score can be written as:

$$\frac{da_e}{d\tau} = \frac{1}{\lambda}(\frac{dh_e}{d\tau})^\top V + \frac{1}{\lambda}h_e^\top \frac{dV}{d\tau}$$

$$= (\frac{de}{d\tau})^\top (W^e)^\top \frac{V}{\lambda} + e^\top (\frac{dW^e}{d\tau})^\top \frac{V}{\lambda} \qquad (20)$$

$$+ \frac{1}{\lambda}h_e^\top \frac{dV}{d\tau}.$$

Let us expand the terms as:

$$
\frac{da_e}{d\tau} = \frac{1}{m} \sum_{(t,j):e_j^{(t)}\equiv e} y^{(t)}\beta^{(t)}\alpha_j^{(t)} \frac{||\boldsymbol{V}^\top \mathbf{W}^e||_2^2}{\lambda^2}(s_e - s)
$$
$$
+ \frac{1}{m} \sum_{(t,j):e_j^{(t)}\equiv e} y^{(t)}\beta^{(t)}\alpha_j^{(t)} \frac{\boldsymbol{W'}^\top \boldsymbol{V'}}{\lambda}
$$
$$
+ \frac{1}{m} \sum_{t=1}^{m} \sum_j y^{(t)}\beta^{(t)}\alpha_j^{(t)} \boldsymbol{e}^\top \boldsymbol{e}_j^{(t)} \frac{\boldsymbol{W}^\top \boldsymbol{V}}{\lambda}
$$
$$
+ \frac{1}{m} \sum_{t=1}^{m} \sum_j y^{(t)}\beta^{(t)}\alpha_j^{(t)}(s_j^{(t)} - s^{(t)})\boldsymbol{e}^\top \boldsymbol{e}_j^{(t)} \frac{||\boldsymbol{V}||_2^2}{\lambda^2}
$$
$$
+ \frac{1}{m} \sum_{t=1}^{m} \sum_j y^{(t)}\beta^{(t)}\alpha_j^{(t)}(s_j^{(t)} - s^{(t)}) \frac{\boldsymbol{h}_e^\top \boldsymbol{h}_j^{(t)}}{\lambda^2},
\tag{21}
$$

where $\boldsymbol{W'} = (\mathbf{W}^e)^\top \boldsymbol{W}$, $\boldsymbol{V'} = (\mathbf{W}^e)^\top \boldsymbol{V}$.

## C  Model with Additive Attention

Additive attention is commonly used for text classifications. Let us replace the scaled dot-product attention in the model discussed in the paper with additive attention. To make it more general, we also consider a scaling factor for the additive attention. The attention score at the $i$-th token of an instance will be calculated as:

$$
a_i = \tanh\left(\boldsymbol{e}_i^\top \mathbf{W}^a + \boldsymbol{b}^a\right)\boldsymbol{V}/\lambda,
\tag{22}
$$

where $\mathbf{W}^a \in R^{d\times d}$ is the weight matrix, $\boldsymbol{b}^a \in R^d$ is the bias for the attention layer. For simplicity, we do not consider the bias here, and the attention score will be:

$$
a_i = \frac{\tanh\left(\boldsymbol{e}_i^\top \mathbf{W}^a\right)\boldsymbol{V}}{\lambda}.
\tag{23}
$$

Let us use $\boldsymbol{x}_i$ to denote $\boldsymbol{e}_i^\top \mathbf{W}^a$, and:

$$
a_i = \frac{\tanh(\boldsymbol{x}_i)\boldsymbol{V}}{\lambda}.
\tag{24}
$$

The derivative of the loss with respect to the embedding $\boldsymbol{e}_i$ will be calculated as:

$$
\frac{\partial \ell}{\partial \boldsymbol{e}_i} = -y\beta\alpha_i \left[\frac{\mathbf{W}^a \mathbf{D}_i \boldsymbol{V}(\boldsymbol{h}_i - \boldsymbol{h})^\top}{\lambda} + \mathbf{I}\right]\boldsymbol{W},
\tag{25}
$$

where $\mathbf{D}_i = \text{Diag}(1 - \tanh^2 \boldsymbol{x}_i)$, which is a diagonal matrix. The derivative of the loss with respect to $\boldsymbol{V}$ will be calculated as:

$$
\frac{\partial \ell}{\partial \boldsymbol{V}} = -\frac{y\beta}{\lambda} \sum_j \alpha_j s_j(\tanh \boldsymbol{x}_j^\top - \sum_i \alpha_i \tanh \boldsymbol{x}_i^\top).
\tag{26}
$$

And the derivative of the loss with respect to $\mathbf{W}^a$ will be:

$$
\frac{\partial \ell}{\partial \mathbf{W}^a} = -\frac{y\beta}{\lambda} \sum_j \alpha_j s_j(\boldsymbol{e}_j \boldsymbol{V}^\top \mathbf{D}_j - \sum_i \alpha_i \boldsymbol{e}_i \boldsymbol{V}^\top \mathbf{D}_i).
\tag{27}
$$

The update of the polarity score can be written as:

$$
\frac{ds_e}{d\tau} = \left(\frac{d\boldsymbol{e}}{d\tau}\right)^\top \boldsymbol{W} + \boldsymbol{e}^\top \frac{d\boldsymbol{W}}{d\tau}.
\tag{28}
$$

When the scaling factor $\lambda$ is large enough, the derivative can be approximated in a way similar to the model with scaled dot-product attention. And we will make similar predictions on the trends of the polarity tokens as well as the neutral tokens.

However, the update of the attention score will be different from the scaled dot-product one as below:

$$
\frac{da_e}{d\tau} = \frac{d\tanh \boldsymbol{x}_e^\top}{d\tau} \frac{\boldsymbol{V}}{\lambda} + \frac{\tanh \boldsymbol{x}_e^\top}{\lambda} \frac{d\boldsymbol{V}}{d\tau}
\tag{29}
$$

It is complex to handle the function $\tanh$ directly.

Alternatively, we can approximate the function $\tanh$ using first-order Taylor series, namely $\tanh(x) \approx x$. The update of attention score on the entire training set can be approximated as:

$$
\frac{da_e}{d\tau} \approx \left(\frac{d\boldsymbol{x}_e}{d\tau}\right)^\top \boldsymbol{V} + \boldsymbol{x}_e^\top \frac{d\boldsymbol{V}}{d\tau}
$$
$$
= \left(\frac{d\boldsymbol{e}}{d\tau}\right)^\top \mathbf{W}^a \frac{\boldsymbol{V}}{\lambda} + \boldsymbol{e}^\top \frac{d\mathbf{W}^a}{d\tau} \frac{\boldsymbol{V}}{\lambda}
$$
$$
+ \frac{1}{\lambda} \boldsymbol{e}^\top \mathbf{W}^a \frac{d\boldsymbol{V}}{d\tau}
$$
$$
\approx \frac{1}{m} \sum_{(t,j):e_j^{(t)}\equiv e} y^{(t)}\beta^{(t)}\alpha_j^{(t)} \frac{||\mathbf{W}^a \boldsymbol{V}||_2^2}{\lambda^2}(s_e - s^{(t)})
$$
$$
+ \frac{1}{m} \sum_{(t,j):e_j^{(t)}\equiv e} y^{(t)}\beta^{(t)}\alpha_j^{(t)} \frac{\boldsymbol{W}^\top \mathbf{W}^a \boldsymbol{V}}{\lambda}
$$
$$
+ \frac{1}{m} \sum_{t=1}^{m} \sum_j y^{(t)}\beta^{(t)}\alpha_j^{(t)}(s_j^{(t)} - s^{(t)})\boldsymbol{e}^\top \boldsymbol{e}_j^{(t)} \frac{||\boldsymbol{V}||_2^2}{\lambda^2}
$$
$$
+ \frac{1}{m} \sum_{t=1}^{m} \sum_j y^{(t)}\beta^{(t)}\alpha_j^{(t)}(s_j^{(t)} - s^{(t)}) \frac{\boldsymbol{x}\boldsymbol{x}_j^\top}{\lambda^2},
\tag{30}
$$

where $\boldsymbol{x} = \boldsymbol{e}^\top \mathbf{W}^a$.

## D  Multi-class Classification

The architecture for multi-class classification is similar to the one discussed in our paper except

that: the last linear layer weight will be a matrix instead of a vector, and the sigmoid activation will be replaced by a softmax function. The token representation is the embedding.

The instance-level polarity scores[1] will be calculated as:

$$\boldsymbol{s} = (\mathbf{W}^h)^\top \boldsymbol{h}, \tag{31}$$

where $\boldsymbol{h}$ is the instance representation, $\mathbf{W}^h \in R^{d \times K}$ is the weight matrix of the final linear layer, $K$ is the class number. $\boldsymbol{s}$ is a vector consisting of $K$ elements, each element will be regarded as the polarity score with respect to its corresponding label. A desirable scenario is such: if the instance is with the label $k$, then the $k$-th element in $\boldsymbol{s}$ will be as large as possible whereas the other elements will be as small as possible .

We can get the probability distribution for all the labels:

$$\boldsymbol{p} = \text{Softmax}(\boldsymbol{s}), \tag{32}$$

where $\boldsymbol{p} \in R^K$.

Cross-entropy loss will be used, and the loss on an instance with the label $k \in (0, 1, .., K-1)$ can be calculated as:

$$\ell = -\log p_k, \tag{33}$$

where $p_k$ refers to the corresponding probability for the label $k$.

The partial derivative of the loss $\ell$ with respect to $\boldsymbol{h}$ will be calculated as:

$$\frac{\partial \ell}{\partial \boldsymbol{h}} = \mathbf{W}^h \hat{\boldsymbol{p}}, \quad \hat{\boldsymbol{p}} = \boldsymbol{p} - \boldsymbol{I}_k, \tag{34}$$

where $\boldsymbol{I}_k \in R^K$ is the one-hot encoding for the $k$-th label. Note that the elements of the vector $\hat{\boldsymbol{p}} \in R^K$ are positive except the one that corresponds to the ground truth label.

Given an instance with the label $k$, we can have such partial derivatives:

$$\frac{\partial \ell}{\partial \boldsymbol{s}} = \hat{\boldsymbol{p}}, \tag{35}$$

$$\frac{\partial \ell}{\partial \mathbf{W}^h} = \boldsymbol{h} \hat{\boldsymbol{p}}^\top. \tag{36}$$

---
[1] For simplicity, we do not consider the bias in the linear layer as we did in the paper.

On the entire dataset, the partial derivatives will be calculated as:

$$\frac{\partial \ell}{\partial \boldsymbol{e}} = \frac{1}{m} \sum_{(t,j): e_j^{(t)} \equiv e} \alpha_j^{(t)} \left[ \frac{\boldsymbol{V}(\boldsymbol{e} - \boldsymbol{h}^{(t)})^\top}{\lambda} + \mathbf{I} \right] \mathbf{W}^h \hat{\boldsymbol{p}}^{(t)}, \tag{37}$$

$$\frac{\partial \ell}{\partial \mathbf{W}^h} = \frac{1}{m} \sum_{t=1}^m \boldsymbol{h}^{(t)} (\hat{\boldsymbol{p}}^{(t)})^\top. \tag{38}$$

On the entire training set, the gradient of $\boldsymbol{V}$ can be calculated as:

$$\frac{\partial \ell}{\partial \boldsymbol{V}} = \frac{1}{m\lambda} \sum_{t=1}^m \sum_j \alpha_j^{(t)} \boldsymbol{e}_j^{(t)} \left( \boldsymbol{e}_j^{(t)} - \boldsymbol{h}^{(t)} \right)^\top \mathbf{W}^h \hat{\boldsymbol{p}}^{(t)}$$

$$= \frac{1}{m\lambda} \sum_{t=1}^m \sum_j \alpha_j^{(t)} \boldsymbol{e}_j^{(t)} \left( \boldsymbol{s}_j^{(t)} - \boldsymbol{s}^{(t)} \right)^\top \hat{\boldsymbol{p}}^{(t)}. \tag{39}$$

Accordingly, given a token $e$, the token-level polarity scores will be calculated as:

$$\boldsymbol{s}_e = (\mathbf{W}^h)^\top \boldsymbol{e}. \tag{40}$$

Note that $\boldsymbol{s}_e$ is a vector that has $K$ elements, each corresponds the polarity score with respect to the label.

The update of $\boldsymbol{s}_e$ will be written as:

$$\frac{d\boldsymbol{s}_e}{d\tau} = \frac{d(\mathbf{W}^h)^\top}{d\tau} \boldsymbol{e} + (\mathbf{W}^h)^\top \frac{d\boldsymbol{e}}{d\tau}$$

$$= -\frac{1}{m\lambda} \sum_{(t,j): e_j^{(t)} \equiv e} (\mathbf{W}^h)^\top \boldsymbol{V}(\boldsymbol{s}_e - \boldsymbol{s}^{(t)})^\top \hat{\boldsymbol{p}}^{(t)}$$

$$- \frac{1}{m} \sum_{(t,j): e_j^{(t)} \equiv e} (\mathbf{W}^h)^\top \mathbf{W}^h \hat{\boldsymbol{p}}^{(t)}$$

$$- \frac{1}{m} \sum_{t=1}^m \hat{\boldsymbol{p}}^{(t)} (\boldsymbol{h}^{(t)})^\top \boldsymbol{e}. \tag{41}$$

And the update of $a_e$ wil be written as:

$$\frac{da_e}{d\tau} = \frac{1}{\lambda} \boldsymbol{e}^\top \frac{d\boldsymbol{V}}{d\tau} + \frac{1}{\lambda} \left( \frac{d\boldsymbol{e}}{d\tau} \right)^\top \boldsymbol{V}$$

$$= \frac{1}{m\lambda} \sum_{t=1}^m \sum_j \alpha_j^{(t)} \boldsymbol{e}^\top \boldsymbol{e}_j^{(t)} \left( \boldsymbol{s}_j^{(t)} - \boldsymbol{s}^{(t)} \right)^\top \hat{\boldsymbol{p}}^{(t)}$$

$$- \frac{1}{m\lambda} \sum_{(t,j): e_j^{(t)} \equiv e} \alpha_j^{(t)} (\hat{\boldsymbol{p}}^{(t)})^\top \frac{(\boldsymbol{s}_e - \boldsymbol{s}^{(t)}) \|\boldsymbol{V}\|_2^2}{\lambda}$$

$$- \frac{1}{m\lambda} \sum_{(t,j): e_j^{(t)} \equiv e} \alpha_j^{(t)} (\hat{\boldsymbol{p}}^{(t)})^\top (\mathbf{W}^h)^\top \boldsymbol{V}. \tag{42}$$