

Intelligent Speech Distinguish Homework 02

Qiu Yihang

June.13-17, 2022

Backward Propagation in RNN.

Input Sequence: \mathbf{x} . Output Sequence: $\hat{\mathbf{r}}$. Label Sequence: \mathbf{r} (with length T_r).

Number of categories: C . Activation function: $\sigma(z) = \frac{1}{1+e^{-z}}$.

Network Structure:

- Input Layer: $\mathbf{a}_t^{(\text{in})} = \sigma(\mathbf{W}^{(\text{in})}\mathbf{x}_t + \mathbf{b}^{(\text{in})})$.
- Hidden Layer (RNN): $\mathbf{h}_t = \sigma(\mathbf{U}\mathbf{a}_t^{(\text{in})} + \mathbf{V}\mathbf{h}_{t-1} + \mathbf{b}_h)$, $\mathbf{o}_t = \sigma(\mathbf{W}\mathbf{h}_t + \mathbf{b}_o)$.
- Output Layer: $\mathbf{h}_t^{(\text{out})} = \mathbf{W}^{(\text{out})}\mathbf{o}_t + \mathbf{b}^{(\text{out})}$, $\hat{\mathbf{r}}_t = \text{softmax}(\mathbf{h}_t^{(\text{out})})$.
- Loss Function: (Cross Entropy Loss)

$$\mathcal{L} = \sum_{t=1}^{T_r} \mathcal{L}_{\text{oss}}(r_t, \hat{\mathbf{r}}_t) = - \sum_{t=1}^{T_r} \sum_{i=1}^C r_{t,i} \log \hat{\mathbf{r}}_{t,i} = - \sum_{t=1}^{T_r} r_t^T \log \hat{\mathbf{r}}_t$$

Give the Backward Propagation of RNN.

Proof. Use $\delta\alpha$ to denote $\frac{\partial \mathcal{L}}{\partial \alpha}$. We know $\sigma'(z) = \sigma(z)(1 - \sigma(z))$.

Then we have

$$\begin{aligned}
 \delta \hat{\mathbf{r}}_{t,i} &= \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{r}}_{t,i}} = \frac{r_{t,i}}{\hat{\mathbf{r}}_{t,i}} \quad (1 \leq i \leq C) \quad \Longleftrightarrow \quad \delta \hat{\mathbf{r}}_t = r_t \oslash \hat{\mathbf{r}}_t \quad (\text{where } \oslash \text{ is element-wise division}) \\
 \delta \mathbf{h}_{t,i}^{(\text{out})} &= \sum_{j=1}^C \delta \hat{\mathbf{r}}_{t,j} \cdot \text{softmax}'_{i,j}(\mathbf{h}_{t,i}^{(\text{out})}) = - \sum_{j \neq i} \frac{r_{t,j}}{\hat{\mathbf{r}}_{t,j}} (-\hat{\mathbf{r}}_{t,i} \hat{\mathbf{r}}_{t,j}) - \frac{r_{t,i}}{\hat{\mathbf{r}}_{t,i}} \hat{\mathbf{r}}_{t,i} (1 - \hat{\mathbf{r}}_{t,i}) \\
 &= -r_{t,i} + \hat{\mathbf{r}}_{t,i} \sum_{j=1}^C r_{t,j} = \hat{\mathbf{r}}_{t,i} - r_{t,i} \quad \Longleftrightarrow \quad \delta \mathbf{h}_t^{(\text{out})} = \hat{\mathbf{r}}_t - r_t \\
 \delta \mathbf{o}_t &= \frac{\partial \mathbf{h}_t^{(\text{out})}}{\partial \mathbf{o}_t} \delta \mathbf{h}_t^{(\text{out})} = \mathbf{W}^{(\text{out})} \delta \mathbf{h}_t^{(\text{out})} \\
 \delta \mathbf{b}^{(\text{out})} &= \delta \mathbf{h}_t^{(\text{out})} \\
 \delta \mathbf{W}_{i,j}^{(\text{out})} &= \mathbf{o}_{t,j} \cdot \delta \mathbf{h}_{t,i}^{(\text{out})} \\
 \delta \mathbf{y}_o &\triangleq \frac{\partial \mathcal{L}}{\partial (\mathbf{W}\mathbf{h}_t + \mathbf{b}_o)} \\
 &= \sigma(\mathbf{W}\mathbf{h}_t + \mathbf{b}_o) \odot (1 - \sigma(\mathbf{W}\mathbf{h}_t + \mathbf{b}_o)) \odot \delta \mathbf{o}_t \quad (\text{where } \odot \text{ is element-wise product})
 \end{aligned}$$

$$\begin{aligned}
\delta \mathbf{h}_t &= \delta \mathbf{y}_o \frac{\partial (\mathbf{W} \mathbf{h}_t + \mathbf{b}_o)}{\partial \mathbf{h}_t} = \mathbf{W} \delta \mathbf{y}_o \\
\delta \mathbf{b}_o &= \delta \mathbf{y}_o \frac{\partial (\mathbf{W} \mathbf{h}_t + \mathbf{b}_o)}{\partial \mathbf{b}_o} = \delta \mathbf{y}_o \\
\delta \mathbf{W}_{i,j} &= \delta \mathbf{y}_{o,i} \frac{\partial (\mathbf{W} \mathbf{h}_t + \mathbf{b}_o)_i}{\partial \mathbf{W}_{i,j}} = \mathbf{h}_{t,j} \cdot \delta \mathbf{y}_{o,i} \\
\delta \mathbf{y}_h &\triangleq \frac{\partial \mathcal{L}}{\partial (\mathbf{U} \mathbf{a}_t^{(\text{in})} + \mathbf{V} \mathbf{h}_{t-1} + \mathbf{b}_h)} \\
&= \sigma (\mathbf{U} \mathbf{a}_t^{(\text{in})} + \mathbf{V} \mathbf{h}_{t-1} + \mathbf{b}_h) \odot (\mathbf{1} - \sigma (\mathbf{U} \mathbf{a}_t^{(\text{in})} + \mathbf{V} \mathbf{h}_{t-1} + \mathbf{b}_h)) \odot \delta \mathbf{h}_t \\
\delta \mathbf{a}_t^{(\text{in})} &= \mathbf{U} \delta \mathbf{y}_h \\
\delta \mathbf{b}_h &= \delta \mathbf{y}_h \\
\delta \mathbf{U}_{i,j} &= \mathbf{a}_{t,j}^{(\text{in})} \cdot \delta \mathbf{y}_{h,i} \\
\delta \mathbf{V}_{i,j} &= \mathbf{h}_{t-1,j} \cdot \delta \mathbf{y}_{h,i} \\
\delta \mathbf{y}_a &\triangleq \frac{\partial \mathcal{L}}{\partial (\mathbf{W}^{(\text{in})} \mathbf{x}_t + \mathbf{b}^{(\text{in})})} = \sigma (\mathbf{W}^{(\text{in})} \mathbf{x}_t + \mathbf{b}^{(\text{in})}) \odot (\mathbf{1} - \sigma (\mathbf{W}^{(\text{in})} \mathbf{x}_t + \mathbf{b}^{(\text{in})})) \odot \delta \mathbf{a}^{(\text{in})} \\
\delta \mathbf{b}^{(\text{in})} &= \delta \mathbf{y}_a \\
\delta \mathbf{W}_{i,j}^{(\text{in})} &= \mathbf{x}_{t,j} \cdot \delta \mathbf{y}_{a,i}
\end{aligned}$$

In conclusion,

Define $\delta \mathbf{h}_t^{(\text{out})} = \hat{\mathbf{r}}_t - r_t$.

Use \mathbf{v}_i to denote the i -th element of vector \mathbf{v} .

Use $\mathbf{A}_{i,j}$ to denote the element in the i -th row and j -th column of matrix \mathbf{A} .

Then the gradient of all parameters trainable during the training process are as follows.

$$\left\{ \begin{array}{ll}
\delta \mathbf{b}^{(\text{out})} = \delta \mathbf{h}_t^{(\text{out})} & \text{(a vector)} \\
\delta \mathbf{W}_{i,j}^{(\text{out})} = \mathbf{o}_{t,j} \cdot (\delta \mathbf{h}_t^{(\text{out})})_i & \text{(a number) for any element } \mathbf{W}_{i,j}^{(\text{out})} \text{ in } \mathbf{W}^{(\text{out})} \\
\delta \mathbf{b}_o = \sigma' (\mathbf{W} \mathbf{h}_t + \mathbf{b}_o) \odot \delta \mathbf{o}_t & \text{(a vector)} \\
\delta \mathbf{W}_{i,j} = (\mathbf{h}_t)_j \cdot (\delta \mathbf{b}_o)_i & \text{(a number) for any element } \mathbf{W}_{i,j} \text{ in } \mathbf{W} \\
\delta \mathbf{b}_h = \sigma' (\mathbf{U} \mathbf{a}_t^{(\text{in})} + \mathbf{V} \mathbf{h}_{t-1} + \mathbf{b}_h) \odot \mathbf{W} \delta \mathbf{b}_o & \text{(a vector)} \\
\delta \mathbf{U}_{i,j} = (\mathbf{a}_t^{(\text{in})})_j \cdot (\delta \mathbf{b}_h)_i & \text{(a number) for any element } \mathbf{U}_{i,j} \text{ in } \mathbf{U} \\
\delta \mathbf{V}_{i,j} = (\mathbf{h}_{t-1})_j \cdot (\delta \mathbf{b}_h)_i & \text{(a number) for any element } \mathbf{V}_{i,j} \text{ in } \mathbf{V} \\
\delta \mathbf{b}^{(\text{in})} = \sigma' (\mathbf{W}^{(\text{in})} \mathbf{x}_t + \mathbf{b}^{(\text{in})}) \odot \mathbf{U} \delta \mathbf{b}_h & \text{(a vector)} \\
\delta \mathbf{W}_{i,j}^{(\text{in})} = (\mathbf{x}_t)_j \cdot (\delta \mathbf{b}^{(\text{in})})_i & \text{(a number) for any element } \mathbf{W}_{i,j}^{(\text{in})} \text{ in } \mathbf{W}^{(\text{in})}
\end{array} \right.$$

(where $\sigma'(\mathbf{z}) \triangleq \sigma(\mathbf{z}) \odot (\mathbf{1} - \sigma(\mathbf{z}))$)

The backward propagation is as follows.

(where η is the learning rate)

$$\left\{ \begin{array}{l} \mathbf{b}^{(\text{out})} \leftarrow \mathbf{b}^{(\text{out})} - \eta \cdot \delta \mathbf{b}^{(\text{out})} \\ \mathbf{W}_{i,j}^{(\text{out})} \leftarrow \mathbf{W}_{i,j}^{(\text{out})} - \eta \cdot \delta \mathbf{W}_{i,j}^{(\text{out})} \\ \mathbf{b}_o \leftarrow \mathbf{b}_o - \eta \cdot \delta \mathbf{b}_o \\ \mathbf{W}_{i,j} \leftarrow \mathbf{W}_{i,j} - \eta \cdot \delta \mathbf{W}_{i,j} \\ \mathbf{b}_h \leftarrow \mathbf{b}_h - \eta \cdot \delta \mathbf{b}_h \\ \mathbf{U}_{i,j} \leftarrow \mathbf{U}_{i,j} - \eta \cdot \delta \mathbf{U}_{i,j} \\ \mathbf{V}_{i,j} \leftarrow \mathbf{V}_{i,j} - \eta \cdot \delta \mathbf{V}_{i,j} \\ \mathbf{b}^{(\text{in})} \leftarrow \mathbf{b}^{(\text{in})} - \eta \cdot \delta \mathbf{b}^{(\text{in})} \\ \mathbf{W}_{i,j}^{(\text{in})} \leftarrow \mathbf{W}_{i,j}^{(\text{in})} - \eta \cdot \delta \mathbf{W}_{i,j}^{(\text{in})} \end{array} \right.$$

■