

作业3=

1- ① $\frac{32-5}{1} + 1 = 28$, 故为 $8 \times 28 \times 28$

② $\frac{28-2}{2} + 1 = 14$, 故为 $8 \times 14 \times 14$

③ $\frac{14-5}{1} + 1 = 10$, 故为 $16 \times 10 \times 10$

④ $\frac{10-2}{2} + 1 = 5$, 故为 $16 \times 5 \times 5$

2) = ① = 5个 ② = 2个

① $q_1 = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} a_1 = \begin{pmatrix} 3 \\ 6 \\ 9 \end{pmatrix} = k_1 = v_1$

② $q_2 = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$
 $(q_1, q_2, q_3, q_4) = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} (a_1, a_2, a_3, a_4) = \begin{pmatrix} 3 & 6 & 9 & 12 \\ 3 & 6 & 9 & 12 \\ 3 & 6 & 9 & 12 \end{pmatrix}$

$(\frac{b}{k})$ 有 $q = k = v$

$\alpha_{11} = \frac{q_1^T k_1}{\sqrt{3}} = \frac{27}{\sqrt{3}} = 9\sqrt{3}$

$\alpha_{12} = \frac{q_1^T k_2}{\sqrt{3}} = \frac{54}{\sqrt{3}} = 18\sqrt{3}$

$\alpha_{13} = \frac{q_1^T k_3}{\sqrt{3}} = \frac{81}{\sqrt{3}} = 27\sqrt{3}$

$\alpha_{14} = \frac{q_1^T k_4}{\sqrt{3}} = \frac{108}{\sqrt{3}} = 36\sqrt{3}$



$$\therefore \sum_i \exp(\alpha_{i,v}) = e^{9\sqrt{3}} + e^{18\sqrt{3}} + e^{27\sqrt{3}} + e^{36\sqrt{3}} = e^{9\sqrt{3}} (1 + e^{9\sqrt{3}} + e^{18\sqrt{3}} + e^{27\sqrt{3}})$$

$$\therefore \hat{\alpha}_{11} = \frac{\alpha_{11}}{\sum} = \frac{1}{1 + e^{9\sqrt{3}} + e^{18\sqrt{3}} + e^{27\sqrt{3}}}$$

$$\hat{\alpha}_{12} = \frac{\alpha_{12}}{\sum} = \frac{e^{9\sqrt{3}}}{1 + e^{9\sqrt{3}} + e^{18\sqrt{3}} + e^{27\sqrt{3}}}$$

$$\hat{\alpha}_{13} = \frac{\alpha_{13}}{\sum} = \frac{e^{18\sqrt{3}}}{1 + e^{9\sqrt{3}} + e^{18\sqrt{3}} + e^{27\sqrt{3}}}$$

$$\hat{\alpha}_{14} = \frac{e^{27\sqrt{3}}}{1 + e^{9\sqrt{3}} + e^{18\sqrt{3}} + e^{27\sqrt{3}}}$$

$$\therefore b_1 = \sum \hat{\alpha}_{1,v} v_i = \frac{1}{1 + e^{9\sqrt{3}} + e^{18\sqrt{3}} + e^{27\sqrt{3}}}$$

$$\begin{pmatrix} 3 + 6e^{9\sqrt{3}} + 9e^{18\sqrt{3}} + 12e^{27\sqrt{3}} \\ 3 + 6e^{9\sqrt{3}} + 9e^{18\sqrt{3}} + 12e^{27\sqrt{3}} \\ 3 + 6e^{9\sqrt{3}} + 9e^{18\sqrt{3}} + 12e^{27\sqrt{3}} \end{pmatrix}$$

3. 1) 通过 LSTM 通过在系统中增加 gates 的方式来控制信息的流量 (增加/删除), 从而使得在相同步的时候不会出现 gradient 增长过快而 explode 或过小而 vanish 的情况。

2) α 大小为 1, 减小了 CNN 模型的参数个数

① 不会丢失图片信息, 即不会像 $\rightarrow 1$ 的 kernel 那样出现图片边缘 pixels 缺失, 即尺寸不变

② 无论再维护长宽信息只需修改卷积核个数即可进行升维/降维。

