

□ Seq2Seq with Attention: Gradient Checklist

□ Обозначения:

- H: Hidden size
- I: Input size = Output_size + Context_size (в декодере)
- C: Context size (2 * encoder hidden size, если BiLSTM)
- A: Attention vector size
- D: Output size (vocab embedding dim)

□ Decoder: Output Layer

- $Y_t = p_{\text{--}} @ W_{\text{output.T}} + B_{\text{output}}$ # $[1 \times D] = [1 \times (H+C)] \times [(H+C) \times D] + [1 \times D]$
- $dY_t = Y^{\wedge}_t - Y_{\text{true}_t}$ # $[1 \times D]$
- $dW_{\text{output}} += dY_t.T @ p_{\text{--}}$ # $[D \times 1] \times [1 \times (H+C)] = [D \times (H+C)]$
- $dB_{\text{output}} += dY_t$ # $[1 \times D]$

□ LayerNorm (на p)

- $p_{\text{--}} = (p - \text{mean}) / \text{sqrt}(\text{var} + \text{eps})$ # Нормализованное p
- $p_{\text{--}} = p_{\text{--}} * \gamma + \beta$ # Выход после LayerNorm
- $d\Gamma += dY_t @ W_{\text{output}} * p_{\text{--}}$ # $[1 \times (H+C)]$
- $dBeta += dY_t @ W_{\text{output}}$ # $[1 \times (H+C)]$
- $dX = (d * \gamma - \text{mean}_d - x^{\wedge} * \text{mean}_d_x^{\wedge}) / \text{std}$ # $[1 \times (H+C)]$

```

□ Decoder LSTM
• dS_t = dX[:H] + o_T attention # [1×H]
• dContext_t = dX[H:] + o_T attention # [1×C]

• dO_t = dS_t * tanh(C_t) * O_t * (1 - O_t)
• dC_t = _dC_t + dS_t * O_t * (1 - tanh(C_t)^2)
• dC~_t = dC_t * I_t * (1 - C~_t^2)
• dI_t = dC_t * C~_t * I_t * (1 - I_t)
• dF_t = dC_t * C_{t-1} * F_t * (1 - F_t)

• dGates = concat(dF, dI, dC~, dO) # [1×4H]
• dW += x_t.T @ dGates # [I×4H]
• dU += h_{t-1}.T @ dGates # [H×4H]
• dB += dGates # [1×4H]

```

```

□ Attention (Dot-product)
•  $\alpha_t = \text{softmax}(\text{score}(s_{t-1}, h_j))$  # [N×1]
•  $\text{context} = \sum \alpha_j * h_j$  # [1×C]

•  $d\text{Context}_t.\text{dot}(h_k) \rightarrow d\alpha_k$  # scalar
•  $dE_{tj} = \sum_k d\alpha_k * \alpha_k * (\delta_{jk} - \alpha_j)$  # softmax grad
•  $dU_{tj} = dE_{tj} * \text{attention\_vector}$  # [1×A]
•  $d\text{Preact} = dU_{tj} * (1 - \tanh^2)$  # [1×A]
•  $dW_e += d\text{Preact}.T @ h_j$  # [A×C]
•  $dW_d += d\text{Preact}.T @ s_{t-1}$  # [A×H]
•  $dV += u_{tj}.T * dE_{tj}$  # [A×1]

```

- Encoder LSTM (forward и backward)
- Тот же процесс, как и в декодере
- $d0_j = dH_j * \tanh(C_j) * 0_j * (1 - 0_j)$
- $dC_j = dH_j * 0_j * (1 - \tanh^2(C_j)) + _dC$
- Остальные – аналогично декодеру

- $dGates_j = \text{concat}(dF, dI, dC, d0)$ # [1×4H]
- $dW_{enc} += x_j.T @ dGates_j$ # [I×4H]
- $dU_{enc} += h_{j-1}.T @ dGates_j$ # [H×4H]
- $dB_{enc} += dGates_j$ # [1×4H]

□ Конец чеклиста