

Deep equilibrium model

Лебедев Михаил Алексеевич

Основная идея Deep sequence model

На вход подается $\mathbf{x}_{1:T} = [x_1, \dots, x_T] \in \mathbb{R}^{T \times p}$

$$\mathbf{z}_{1:T}^{[i+1]} = f_{\theta}^{[i]}(\mathbf{z}_{1:T}^{[i]}; \mathbf{x}_{1:T}) \quad \text{for } i = 0, 1, 2, \dots, L - 1$$

T - длина последовательности

Sequence model - модель $G(\mathbf{x}_{1:T}) = \mathbf{y}_{1:T} \in \mathbb{R}^{T \times \tilde{q}}$

Выходные данные y зависят только от первых x_1, \dots, x_t и не зависят от последующих x_{t+1}, \dots, x_T .

Проблема

При обучении глубокой сети нам нужно много памяти для backpropagation. Но мы знаем, что сети сходятся. Нам проще найти точку эквilibриума, хранить в памяти только ее.

$$\lim_{i \rightarrow \infty} \mathbf{z}_{1:T}^{[i]} = \lim_{i \rightarrow \infty} f_{\theta}(\mathbf{z}_{1:T}^{[i]}; \mathbf{x}_{1:T}) \equiv f_{\theta}(\mathbf{z}_{1:T}^*; \mathbf{x}_{1:T}) = \mathbf{z}_{1:T}^*$$

Forward pass

В отличие от сверточных сетей, мы используем точку эквilibриума вместо значения активации. Для того чтобы нам ее найти, мы решаем следующее уравнение с помощью квазиньютоновских методов.

$$g_{\theta}(\mathbf{z}_{1:T}^*; \mathbf{x}_{1:T}) = \tilde{f}_{\theta}(\mathbf{z}_{1:T}^*; \mathbf{x}_{1:T}) - \mathbf{z}_{1:T}^* \rightarrow 0.$$

$$\mathbf{z}_{1:T}^{[i+1]} = \mathbf{z}_{1:T}^{[i]} - \alpha B g_{\theta}(\mathbf{z}_{1:T}^{[i]}; \mathbf{x}_{1:T}) \quad \text{for } i = 0, 1, 2, \dots$$

Backward pass

Здесь мы используем следующую теорему

Theorem 1. (Gradient of the Equilibrium Model) Let $\mathbf{z}_{1:T}^* \in \mathbb{R}^{T \times d}$ be an equilibrium hidden sequence with length T and dimensionality d , and $\mathbf{y}_{1:T} \in \mathbb{R}^{T \times q}$ the ground-truth (target) sequence. Let $h : \mathbb{R}^d \rightarrow \mathbb{R}^q$ be any differentiable function and let $\mathcal{L} : \mathbb{R}^q \times \mathbb{R}^q \rightarrow \mathbb{R}$ be a loss function (where h, \mathcal{L} are applied in a vectorized manner) that computes

$$\ell = \mathcal{L}(h(\mathbf{z}_{1:T}^*), \mathbf{y}_{1:T}) = \mathcal{L}(h(\text{RootFind}(g_\theta; \mathbf{x}_{1:T})), \mathbf{y}_{1:T}). \quad (7)$$

Then the loss gradient w.r.t. (\cdot) (for instance, θ or $\mathbf{x}_{1:T}$) is

$$\frac{\partial \ell}{\partial (\cdot)} = -\frac{\partial \ell}{\partial \mathbf{z}_{1:T}^*} (J_{g_\theta}^{-1} \big|_{\mathbf{z}_{1:T}^*}) \frac{\partial f_\theta(\mathbf{z}_{1:T}^*; \mathbf{x}_{1:T})}{\partial (\cdot)} = -\frac{\partial \ell}{\partial h} \frac{\partial h}{\partial \mathbf{z}_{1:T}^*} (J_{g_\theta}^{-1} \big|_{\mathbf{z}_{1:T}^*}) \frac{\partial f_\theta(\mathbf{z}_{1:T}^*; \mathbf{x}_{1:T})}{\partial (\cdot)}, \quad (8)$$

where $J_{g_\theta}^{-1} \big|_{\mathbf{x}}$ is the inverse Jacobian of g_θ evaluated at \mathbf{x} .

The proof is provided in Appendix [A](#). The insight provided by Theorem [1](#) is at the core of our method and its various benefits. Importantly, the backward gradient through the “infinite” stacking can be represented as one step of matrix multiplication that involves the Jacobian at equilibrium. For instance, an SGD update step on model parameters θ would be

$$\theta^+ = \theta - \alpha \cdot \frac{\partial \ell}{\partial \theta} = \theta + \alpha \frac{\partial \ell}{\partial \mathbf{z}_{1:T}^*} (J_{g_\theta}^{-1} \big|_{\mathbf{z}_{1:T}^*}) \frac{\partial f_\theta(\mathbf{z}_{1:T}^*; \mathbf{x}_{1:T})}{\partial \theta}. \quad (9)$$

Далее мы решаем уравнение для нахождения произведения Якобиана на вектор(решаем с помощью метода Бroyдена)

$$J_{g_\theta}^{-1} \big|_{\mathbf{z}_{1:T}^{[i+1]}} \approx B_{g_\theta}^{[i+1]} = B_{g_\theta}^{[i]} + \frac{\Delta \mathbf{z}^{[i+1]} - B_{g_\theta}^{[i]} \Delta g_\theta^{[i+1]}}{\Delta \mathbf{z}^{[i+1] \top} B_{g_\theta}^{[i]} \Delta g_\theta^{[i+1]}} \Delta \mathbf{z}^{[i+1] \top} B_{g_\theta}^{[i]},$$

Свойства

Эффективное использование памяти - для deer equilibrium model используется константное количество памяти(не зависит от глубины сети)

Выбор функции для backforward - любая ограниченная и гладкая функция.

Trellis network

Trellis network - temporal convolutional network, которая сочетает в себе свойства рекуррентных и сверточных нейронных сетей.

$$\tilde{\mathbf{x}}_{1:T} = \text{Input injection (i.e., linearly transformed inputs by Conv1D}(\mathbf{x}_{1:T}; W_x))$$
$$f_{\theta}(\mathbf{z}_{1:T}; \mathbf{x}_{1:T}) = \psi(\text{Conv1D}([\mathbf{u}_{-(k-1)s:}, \mathbf{z}_{1:T}]; W_z) + \tilde{\mathbf{x}}_{1:T})$$

Weight-tied transformer

Multi-head self-attention transformer со следующими свойствами

- 1) производится регуляризация данных для увеличения обобщающей способности**
- 2) размер модели меньше**
- 3) сеть можно развернуть на любую глубину**

$\tilde{\mathbf{x}}_{1:T}$ = Input injection (i.e., linearly transformed inputs by $\mathbf{x}_{1:T}W_x$)

$$f_{\theta}(\mathbf{z}_{1:T}; \mathbf{x}_{1:T}) = \text{LN}(\phi(\text{LN}(\text{SelfAttention}(\mathbf{z}_{1:T}W_{QKV} + \tilde{\mathbf{x}}_{1:T}; \text{PE}_{1:T}))))$$

Источники

<https://arxiv.org/pdf/1909.01377.pdf>