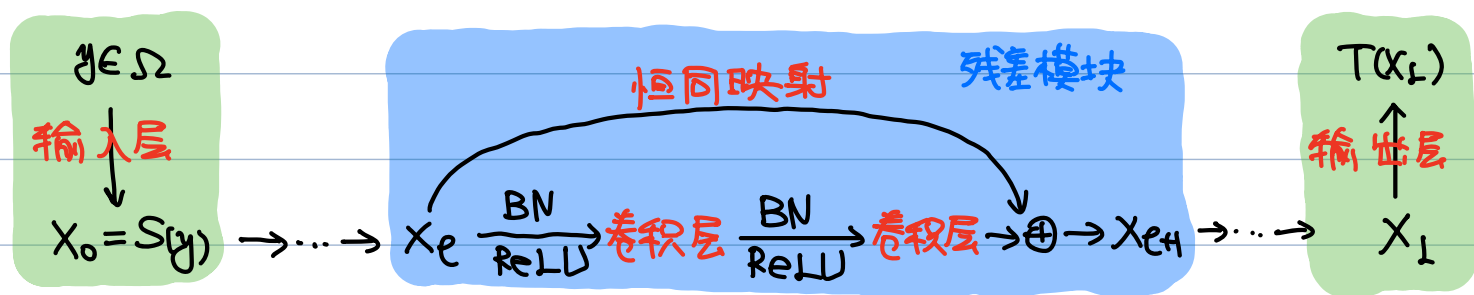


## 1. 残差网络

残差网络 (ResNet) 是深度学习中里程碑式的模型架构。不妨考虑分类学习问题, 即给定输入数据及其标签  $\{y, h(y)\}_{y \in \Omega}$ , 残差网络的前向传播过程如下:



在得到神经网络对输入数据  $y \in \Omega$  的分类预测  $T(X_L)$  后, 便可将其与真实标签  $h(y)$  进行比较, 从而得到损失函数 (例如交叉熵)

$$\mathcal{L}(X_L) = \mathbb{E}_{y \in \Omega} [\|T(X_L) - h(y)\|]$$

不同于卷积层的简单堆叠, 残差模块中引入了恒同映射

$$X_{e+1} = X_e + F(X_e, W_e), \quad 0 \leq e \leq L-1$$

其中  $W_e$  代表模型参数, 从而有效解决了梯度爆炸的问题。

综上所述, 用残差网络做分类任务的优化问题可写为

$$\argmin_{\{W_e\}_{e=0}^L} \left\{ \mathcal{L}(X_L) \mid X_0 = S(y); X_{e+1} = X_e + F(X_e, W_e), 0 \leq e \leq L-1 \right\}$$

其中模型参数的更新采用经典的反向传播算法, 即

$$W_e \leftarrow W_e - \eta \frac{\partial \mathcal{L}(X_L)}{\partial X_{e+1}} \frac{\partial X_{e+1}}{\partial W_e}, \quad 0 \leq e \leq L-1.$$

其中  $\eta > 0$  代表学习率。

## 2. 神经微分方程

若将残差模块看作是常微分方程的前向 Euler 离散格式, 则其相应的优化问题可写为:

$$\operatorname{argmin}_{w(t)} \left\{ \varphi(x(t_1)) \mid x(0) = S_{\varphi}; \frac{dx(t)}{dt} = f(t, x(t), w(t)), 0 < t \leq 1 \right\}$$

在科学计算领域, 上述问题可写为更一般的形式

$$\begin{aligned} & \text{minimize } \overset{\text{损失泛函}}{J[w(t)]} = \varphi(x(t_1), t_1) + \int_{t_0}^{t_1} L(t, x(t), w(t)) dt \\ & \text{subject to } x(t_0) = x_0 \in \mathbb{R}^{m \times 1}; \quad \frac{dx(t)}{dt} = f(t, \underset{\text{状态变量}}{x(t)}, \underset{\text{控制变量}}{w(t)}), t_0 < t \leq t_1 \end{aligned}$$

为了求解上述带约束的最优化问题, 引入 Lagrange 乘子

$$\text{伴随变量} \leftarrow p(t) = [p_1(t), p_2(t), \dots, p_m(t)]^T \in \mathbb{R}^{1 \times m}$$

将其改写为无约束的最优化问题, 即增广 Lagrange 泛函

$$\begin{aligned} J_a[x, w, p] &= \varphi(x(t_1), t_1) + \int_{t_0}^{t_1} [L(t, x, w) + p(f(t, x, w) - \dot{x})] dt \\ & \quad (\text{为了符号简便, 省略记号 } "t)", \text{ 并将对 } t \text{ 求导记为 } "\dot{\cdot}") \end{aligned}$$

$$= \varphi(x(t_1), t_1) + \int_{t_0}^{t_1} [L + p f] dt - \int_{t_0}^{t_1} p \dot{x} dt$$

$$\text{分部积分: } \int_{t_0}^{t_1} p \dot{x} dt = \int_{t_0}^{t_1} p dx = p x \Big|_{t_0}^{t_1} - \int_{t_0}^{t_1} \dot{p} x dt$$

$$= \varphi(x(t_1), t_1) + \int_{t_0}^{t_1} [L + p f + \dot{p} x] dt - p(t_1) x(t_1) + p(t_0) x(t_0)$$

记微分方程中控制变量微小扰动  $\delta w$  造成状态变量的变化为  $\delta x$ , 则对应增广 Lagrange 泛函的变化为

$$\delta J_a = \left[ \left( \frac{\partial \varphi}{\partial x} - p \right) \delta x \right]_{t=t_1} + \int_{t_0}^{t_1} \left[ \left( \frac{\partial L}{\partial x} + p \frac{\partial f}{\partial x} + \dot{p} \right) \delta x + \left( \frac{\partial L}{\partial w} + p \frac{\partial f}{\partial w} \right) \delta w \right] dt$$

其中  $x(t_0)$  为固定的, 因此  $\delta x|_{t=t_0} = 0$ . 若 Lagrange 乘子满足

$$\begin{cases} \mathcal{J}(t_1) = \frac{\partial \mathcal{J}(x(t_1), t_1)}{\partial x}, \\ \frac{d\mathcal{J}(t)}{dt} = - \frac{\partial L(t, x(t), w(t))}{\partial x} - \mathcal{J}(t) \frac{\partial f(t, x(t), w(t))}{\partial x}, \quad t_0 \leq t < t_1, \end{cases}$$

则增广 Lagrange 泛函的变化为

$$\delta J_a = \int_{t_0}^{t_1} \left( \frac{\partial L}{\partial w} + \mathcal{J} \frac{\partial f}{\partial w} \right) \delta w \, dt$$

因此  $w^*(t)$  使得增广 Lagrange 泛函取得最小的必要条件为

$$\left[ \frac{\partial L(t, x(t), w(t))}{\partial w} + \mathcal{J}(t) \frac{\partial f(t, x(t), w(t))}{\partial w} \right]_{w(t)=w^*(t)} = 0, \quad t_0 \leq t \leq t_1.$$

## 神经微分方程的最优控制问题

特别的, 考虑  $\mathcal{J}(x(t_1), t_1) = \mathcal{G}(x(t_1))$ ,  $L(t, x(t), w(t)) = 0$ . 则神经微分方程最优控制问题的解满足,

$$\begin{cases} x^*(0) = x_0; \quad \frac{dx^*(t)}{dt} = f(t, x^*(t), w^*(t)), \quad 0 \leq t \leq 1 & \text{[状态方程]} \\ \mathcal{J}^*(1) = \frac{d\mathcal{G}(x^*(1))}{dx}; \quad \frac{d\mathcal{J}^*(t)}{dt} = -\mathcal{J}^*(t) \frac{\partial f(t, x^*(t), w^*(t))}{\partial x}, \quad 0 \leq t < 1 & \text{[伴随方程]} \\ \mathcal{J}^*(t) \frac{\partial f(t, x^*(t), w^*(t))}{\partial w} = 0, \quad 0 \leq t \leq 1 & \text{[控制方程]} \end{cases}$$

若改写为迭代格式, 即为

$$\begin{cases} x(0) = x_0; \quad dx(t) = f(t, x(t), w(t)), \quad 0 \leq t \leq 1 & \text{[前向传播]} \\ \mathcal{J}(1) = \frac{d\mathcal{G}(x(1))}{dx}; \quad d\mathcal{J}(t) = -\mathcal{J}(t) \frac{\partial f(t, x(t), w(t))}{\partial x}, \quad 0 \leq t < 1 \\ w(t) \leftarrow w(t) - \eta \mathcal{J}(t) \frac{\partial f(t, x(t), w(t))}{\partial w}, \quad 0 \leq t \leq 1 \end{cases} \quad \text{[反向传播]}$$

这与离散情形下 ResNet 的前向-反向传播完全一致,

### 3、应用与扩展

自然而然的, 可以用神经微分方程来解决分类问题、生成模型等诸多应用问题, 有着广泛的应用场景。特别的, 由神经网络表达的右端项  $f(t, x(t), w(t))$  拥有强大的拟合能力, 可根据学习任务演化非常复杂的动力学行为 (详情参见 Jupyter Notebook)。

当然, 神经微分方程方法也有局限性, 需要进一步改进以增强其表达能力, 例如 Augmented Neural ODEs (详情参见 Jupyter Notebook)。