

META-LEARNING WITH NEGATIVE LEARNING RATES

沈雯杰

LAMDA, Nanjing University

January 29, 2021



目录

背景

MIXED LINEAR REGRESSION

TWO CASES

目录

背景

MIXED LINEAR REGRESSION

TWO CASES

MAML

We assume the existence of a distribution of tasks τ and, for each task, a distribution of data points \mathcal{D}^r and a loss function \mathcal{L}^τ . The loss function of the meta-learning problem, $\mathcal{L}^{\text{meta}}$, is defined as an average over both distributions of tasks and data points. The goal of meta-learning is to minimize the loss function with respect to a vector of meta-parameters ω

$$\mathcal{L}^{\text{meta}}(\omega) = \mathbb{E}_{\tau} \mathbb{E}_{\mathcal{D}_t^r \mathcal{D}_v^\tau} \mathbb{E} \mathcal{L}^\tau(\boldsymbol{\theta}^\tau(\omega, \mathcal{D}_t^\tau); \mathcal{D}_v^\tau)$$

$$\mathcal{L}^{\text{meta}}(\omega) = \frac{1}{mn_v} \sum_{i=1}^m \sum_{j=1}^{n_v} \mathcal{L}\left(\boldsymbol{\theta}^{(i)}(\omega); \mathcal{D}_j^{(i)}\right)$$

MAML

In this work we consider the simple case of a single gradient step. Therefore, the inner loop of MAML is given by

$$\theta^{(i)}(\omega) = \omega - \frac{\alpha_t}{n_t} \sum_{j=1}^{n_t} \left. \frac{\partial \mathcal{L}^{(i)}}{\partial \theta} \right|_{\omega; \mathcal{D}_j^{(i)}}$$

目录

背景

MIXED LINEAR REGRESSION

TWO CASES

MIXED LINEAR REGRESSION

In mixed linear regression, each task is characterized by a different linear function, and a model is evaluated by the mean squared error loss function. We assume a generative model in the form of $y = \mathbf{x}^T \mathbf{w} + z$, where \mathbf{x} is the input vector (of dimension p), y is the output (scalar), z is noise (scalar), and \mathbf{w} is a vector of generating parameters (of dimension p)

$$\mathbf{w} \sim \mathcal{N}\left(\mathbf{w}_0, \frac{\nu^2}{p} I_p\right) \quad \mathbf{x} \sim \mathcal{N}(0, I_p) \quad y | \mathbf{x}, \mathbf{w} \sim \mathcal{N}(\mathbf{x}^T \mathbf{w}, \sigma^2)$$

目录

背景

MIXED LINEAR REGRESSION

TWO CASES

OVERPARAMETERIZED CASE

Let $p > n_v m$.

$$\begin{aligned}\bar{\mathcal{L}}^{\text{test}} = & \frac{\sigma^2}{2} \left(1 + \frac{\alpha_r^2 p}{n_r} \right) + \\ & + h^r \left[\frac{\nu^2}{2} \left(1 + \frac{n_v m}{p} \right) + \frac{1}{2} \left(1 - \frac{n_v m}{p} \right) |\omega_0 - \mathbf{w}_0|^2 + \frac{\sigma^2 n_v m}{2p} \frac{1 + \frac{\alpha_t^2 p}{n_t}}{h^t} \right] + o(\xi^{-3/2})\end{aligned}$$

where we define the following expressions

$$h^t = (1 - \alpha_t)^2 + \alpha_t^2 \frac{p + 1}{n_t}$$

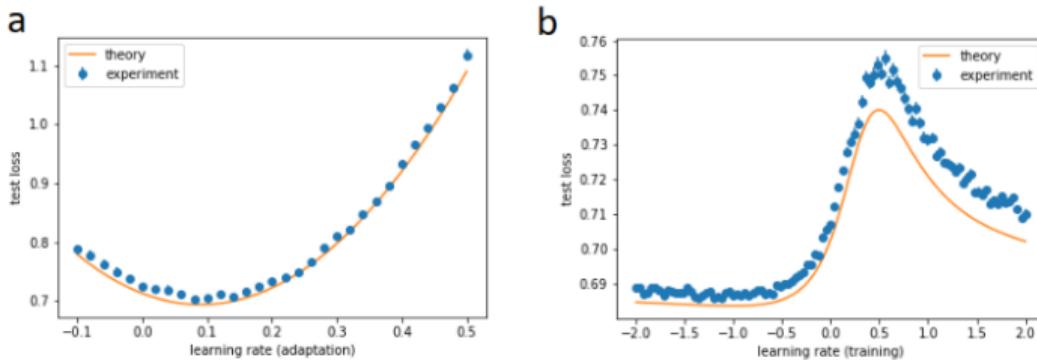
$$h^r = (1 - \alpha_r)^2 + \alpha_r^2 \frac{p + 1}{n_r}$$

Proof. The proof of this Theorem can be found in the Appendix, sections 7.3, 7.3.1.

对于 α_r 来说，这个式子可以看成是两个二次函数的相加，一个最小值取在 0 的时候，一个取在 $\alpha_r = 1 / (1 + (p + 1) / n_r)$ 的时候，所以总的最小值在两者之间并且一定是正的。

对于 α_t 来说，通过求导可以发现它可能在负值处有最小值：

$$\alpha_t^\pm = -\frac{n_t+1}{2p} \pm \sqrt{\left(\frac{n_t+1}{2p}\right)^2 + \frac{n_t}{p}}$$



图：

In panel a) we set $\alpha_t = 0.2$, in panel b) we set $\alpha_r = 0.2$. In the experiments, each run is evaluated on 100 tasks of 50 data points each, and each point is an average over 100 runs (a) or 1000 runs (b).

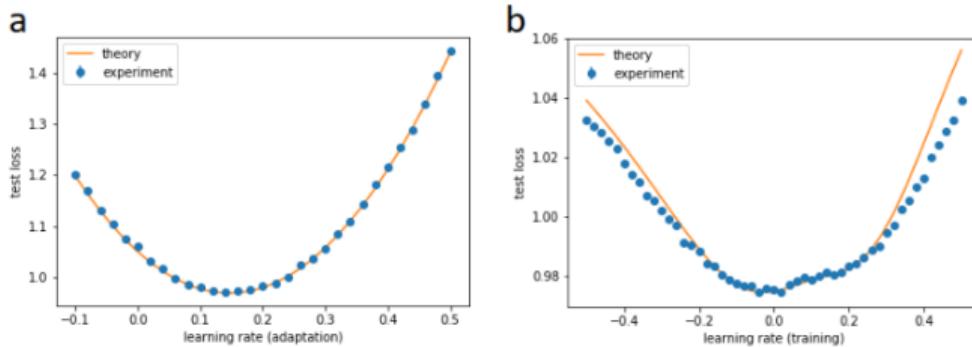
UNDERPARAMETERIZED CASE

Theorem 2. Let $p < n_v m$.

$$\begin{aligned}\bar{\mathcal{L}}^{\text{test}} = & \frac{\sigma^2}{2} \left(1 + \frac{\alpha_r^2 p}{n_r} \right) + \frac{h^r \nu^2}{2} + \\ & + \frac{h^r}{2h^t} \frac{p}{n_v m} \left\{ \sigma^2 \left[h^t + \frac{\alpha_t^2}{n_t} [(n_v + 1) g_1 + pg_2] \right] + \frac{\nu^2}{p} [(n_v + 1) g_3 + pg_4] \right\} + o((m\xi)^{-3/2})\end{aligned}$$

同样的，对于 α_r 来说，这个式子可以看成是两个二次函数的相加，一个最小值取在 0 的时候，一个取在 $\alpha_r = 1 / (1 + (p + 1) / n_r)$ 的时候，所以总的最小值在两者之间并且一定是正的。

对于 α_t 来说，虽然不能直接求出导数为 0 的情况，但是可以得到
$$\frac{\partial \bar{\mathcal{L}}^{\text{test}}}{\partial \alpha_t} \Big|_{\alpha_t=0} = \frac{\sigma^2 p}{n_v m} \geq 0$$
 由此说明可能在负值处取到最小值。



冬

Values of parameters: $n_t = 5$, $n_v = 25$, $n_r = 10$, $m = 40$, $p = 30$, $\sigma = 0.2$, $\nu = 0.2$. In panel a) we set $\alpha_t = 0.2$, in panel b) we set $\alpha_r = 0.2$. In the experiments, the model is evaluated on 100 tasks of 50 data points each, and each point is an average over 100 (a) or 1000 (b) runs.



参考文献

<https://openreview.net/pdf?id=60j5LygnmD>