

3 个训练的超参数的理解

在训练代码里面经常见到的 3 个超参数：

```
1 learning_rate: float = 6e-4
2 weight_decay: float = 1e-1
3 grad_clip: float = 1.0
```

1. learning_rate

2. 为何需要 grad_clip

更新权重的 delta 为： $\text{learning_rate} * \text{grad}$

但是上面的 delta 还是会很大很小，所以继续限制 grad

delta 为： $\text{learning_rate} * \text{Norm_clip}(\text{grad})$

另外：

“The problem is that gradient clipping interacts with adaptive methods like Adam. Since Adam is scale invariant, the absolute scale cannot matter. Rather gradient clipping changes something about the inter-batch variance.”

[Relationship between Learning Rate and Gradient Clipping]

(https://www.reddit.com/r/MachineLearning/comments/eea88q/d_relationship_between_learning_rate_and_gradient/)

Norm_clip:

```
1 if g >= threshold then
2   g <- g*threshold/Norm(g)
```

3. 什么是 weight decay:

$$\text{Loss} = \text{MSE}(\hat{y}, y) + wd * \sum(w^2)$$