

CS224n assignment 3 (Written)

I - (a) - (i)

① How using m stops the updates from varying as much?

: Adam is an incorporated form of Momentum and RMSProp.

The term m is a modified term of velocity. The update rule of momentum (below) contains accumulating the gradients. Not only making steps with respect to current gradients, it also makes steps with respect to previous gradient history. Considering both current and previous gradients, step size will be larger if there are successive gradients with same directions, while smaller with different directions — making variance low.

$$m \leftarrow \beta_1 m + (1 - \beta_1) \nabla_{\theta} J_{\text{minibatch}}(\theta)$$

$$\theta \leftarrow \theta - \alpha m$$

② Why low variance may be helpful to learning?

: Because it helps fast convergence toward loss minima.

I - (a) - (ii)

① Since Adam divides the update by \sqrt{v} , which of the model parameters will get larger updates?

: \sqrt{v} means the squared root sum of all their historical squared values. And according to the update rule below, step size will be larger if \sqrt{v} gets smaller. Small \sqrt{v} means small sum of all historical squared gradients. Therefore, model parameters with small partial derivatives will get larger updates.

$$m \leftarrow \beta_1 m + (1 - \beta_1) \nabla_{\theta} J_{\text{minibatch}}(\theta)$$

$$v \leftarrow \beta_2 v + (1 - \beta_2) (\nabla_{\theta} J_{\text{minibatch}}(\theta) \odot \nabla_{\theta} J_{\text{minibatch}}(\theta))$$

$$\theta \leftarrow \theta - \alpha \odot m / \sqrt{v}$$

② Why might this help learning?

: It can help solving the saddle-point/local-minima problems.

* Intuition!

평균값이 커지 Loss에 민감한

parameter는 천천히 학습되도록 설정

(vice versa)

I - (b)

(i) What must r equal in terms of P_{drop} ?

The value r should have the value that can make the expected value $P_{\text{drop}} \cdot h$.

Therefore, we can compute the expected value as below:

$$E_{\text{Pdrop}}[h_{\text{drop}}]_i = r \cdot 0 \cdot h_i P_{\text{drop}} + r \cdot 1 \cdot h_i (1 - P_{\text{drop}})$$

$$= r \cdot h_i (1 - P_{\text{drop}}) = h_i$$

In conclusion, the value r must be

$$r = \frac{1}{1 - P_{\text{drop}}} \quad \blacksquare$$

(ii) Why should we apply dropout during training but not during evaluation?

: Applying dropout during training means we're ensembling each specific network (including dropout). On evaluation process, we check the ensembled result, so we do not apply dropout during evaluation.