

Lecture 24: Deep neural networks

Professor Ilias Bilonis

The stochastic gradient decent and variants

The Robbins-Monro algorithm

Stochastic optimization algorithm

$$\min_{\theta} \mathbb{E}_Z[\ell(\theta; Z)] (*)$$

1. Initialize θ_0 .

2. Iterate :

$$\theta_{t+1} = \theta_t - \alpha_t \nabla \ell(\theta_t; Z_t)$$

learning rate (scalar)

where Z_t is a sample of Z . (i.i.d.)

RM theorem the algorithm above will converge to a local minimum of the st. opt. (*) if

$$\sum_{t=1}^{\infty} \alpha_t = +\infty ; \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

↳ (α_t goes to zero but not too fast)

$$\alpha_t = \frac{A}{(Bt+C)^p}, \quad 0.5 < p < 1$$

Application of Robbins-Monro to loss minimization

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i; \theta))^2$$

m (batch size), $I_1, \dots, I_m \sim \text{Categorical}(\frac{1}{n}, \dots, \frac{1}{n})$ i.i.d.

$$l_m(\theta; I_{1:m}) = \frac{1}{m} \sum_{j=1}^m (y_{I_j} - f(x_{I_j}; \theta))^2$$

converges to local min of this

$$\min_{\theta} L(\theta) \Leftrightarrow \min_{\theta} \mathbb{E}_{I_{1:m}} [l_m(\theta; I_{1:m})]$$

R.M. \rightarrow - Sample indices i_1, \dots, i_m from I_1, \dots, I_m

$$\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} l_m(\theta_t; i_{1:m})$$

$$\rightarrow \theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} \sum_{j=1}^m \{ y_{i_j} - f(x_{i_j}; \theta_t) \}^2$$

$$= \theta_t - \alpha_t \sum_{j=1}^m 2 (y_{i_j} - f(x_{i_j}; \theta_t)) \cdot \nabla_{\theta} f(x_{i_j}; \theta_t)$$

error
of prediction

gradient of
the model wrt
to all parameters
same dimension as
the parameters

still have
to define
the right
loss function

~~Backprop
TensorFlow
Pytorch~~

Advanced variants of stochastic gradient descent

↳ to avoid having to select the learning rate

- Stochastic gradient descent with momentum.
- AdaGrad.
- Adam (adaptive moment estimation).