

Learning rate decay
step size ↓ gradually

$$\alpha = \frac{1}{1 + \text{decay rate} \times \text{epoch num}} \times \alpha_0$$

$$\alpha = \frac{\text{epoch num}}{2^{\text{epoch num}}} \times \alpha_0$$

exponential decay

$$\alpha = \frac{\kappa}{\sqrt{\text{epoch num}}} \times \alpha_0$$

Discrete decay function
Manual decay

Local Optima

0 gradient

local minima

shallow saddle

high dimension space

problem of plateaus

Batch Normalization

Additional hyperparameter → robust network

→ speed ↑, bigger network

Normalization of input in between layers

① all ϵ - for each layer (mini batch)

② $2^{(p)}$ normal $\times \gamma + \beta$

test data

→ exponentially weighted average over mini batch

not applicable for learning parameters

learning on shifted distⁿ / Covariate shift
~ decoupling between layers

→ robust

add some regularization effect

mini batch - ϵ & ϵ^2 - randomness

Hyperparameter tuning

parameters

- uniform random number
- appropriate scale - log scale
- power - random

Grid vs random

assume one constant
more diverse initialization
coarse to fine

Modeling

Boostrapping one model
training models - parallel

Deep learning framework

Binary
multiclass } classification

Softmax Activation function $t = e^{(z^{L+1})}$

Several linear decision boundary

Hard max → (max → 1 else 0)

loss function

$$L(y, \hat{y}) = - \sum_{j=1}^C y_j \log \hat{y}_j$$

Software framework

1. ease
2. speed
3. freely open