陈艺潇 YiqiaoChen 2024/10/20
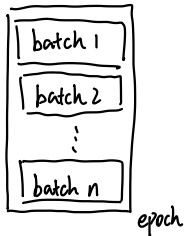
# ① Training Phase

## I. DataSet



batch 1
batch 2
⋮
batch n

epoch

epoch-by-epoch

the entire training process may repeat the epoch thousands or millions of time.

a) Epoch:
⇓
one complete presentation of the entire training set during the learning process.

Randomizing the order of training samples from one epoch to the next tends to avoid the oscillation
( and NN may learn the order of the data samples )

## b) Sequential Learning

Weight updating is performed after the presentation of each training samples

**Advantages**

i. Less Memory Requirements:
Sequential learning don't need to accumulate errors $\{\mathcal{E}_{av}\}_n$ to update like batch learning

ii. Redundancy Utilization:
Each training samples provides a gradient to update the weight, which is batch·size times more frequent than batch learning

iii. Online Stochastic Learning

**Disadvantages**

i. No Theoretical Convergence:
The convergence derivation is all based on batch learning

ii. Limited Parallelization

iii. Sensitivity to Data Order:
NN might learn specific patterns of data order, rather than the generalized feature of the data itself.

## C) Batch Learning

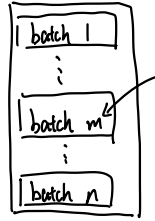Weights are update after the presentation of *all training samples* in a set which is named as batch.

Derivation: $\dfrac{\partial \mathcal{E}_{av}}{\partial e_j} \cdot \dfrac{\partial e_j}{\partial w_{ji}}$

dim = $|x|$

$= \left[ \dfrac{\partial \mathcal{E}_{av}}{\partial e_j(1)} \quad \dfrac{\partial \mathcal{E}_{av}}{\partial e_j(2)} \quad \cdots \quad \dfrac{\partial \mathcal{E}_{av}}{\partial e_j(N)} \right] \begin{bmatrix} \dfrac{\partial e_j(1)}{\partial w_{ji}} \\ \dfrac{\partial e_j(2)}{\partial w_{ji}} \\ \vdots \\ \dfrac{\partial e_j(N)}{\partial w_{ji}} \end{bmatrix}$

$= \sum_{i=1}^{N} \dfrac{\partial \mathcal{E}_{av}}{\partial e_j(i)} \dfrac{\partial e_j(i)}{\partial w_{ji}}$

$= \sum_{i=1}^{N} \dfrac{e_j(i)}{N} \cdot \dfrac{\partial e_j(i)}{\partial w_{ji}}$

| batch 1 |
| ... |
| batch m |
| ... |
| batch n |

in m'th update

$\mathcal{E}_{av} = \dfrac{1}{N} \sum_{i=1}^{N} \mathcal{E}(i) = \dfrac{1}{2N} \sum_{i=1}^{N} \sum_{j \in C} e_j^2(i)$

$\Delta w_{ji} = -\eta \dfrac{\partial \mathcal{E}_{av}}{\partial w_{ji}}$

$= -\dfrac{\eta}{N} \sum_{i=1}^{N} e_j(n) \dfrac{\partial e_j(n)}{\partial w_{ji}}$

### Advantages

i. Accurate Gradients for Convergence:
   The gradient used for batch update is the average of the gradients, which is more *statistically representative*

ii. Parallelization

iii. Improved regularization

### Disadvantages

i. Limited Redundancy Utilization

ii. Memory Requirements

⭐ iii. *Initialization Sensitivity*

iv. Stalling in Local Minima

## II. Stopping Criteria

GD is considered to have converged when

① $\| \nabla_w \mathcal{E}_{av} \|$ reaches a sufficiently small gradient threshold

② $| \Delta \mathcal{E}_{av} |$ is sufficiently small ( 0.1% to 1% )
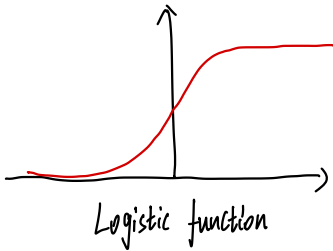
# ② Training Improvement

## I. Stochastic versus batch update

## II. Maximizing information content. Use training samples that
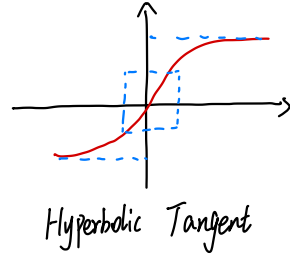  1) result in the largest training error
  2) are radically different from all those previously used  ← distinct samples

## III. Activation function



### Logistic function

i. $\psi(x) = \dfrac{1}{1 + e^{-ax}}$ , $a > 0$

ii. $\psi'(x) = a\,\psi(x)(1 - \psi(x))$

### Hyperbolic Tangent

i. $\psi(x) = a \tanh(bx)$  , $(a, b) > 0$

ii. $\psi'(x) = \dfrac{b}{a}(a - \psi(x))(a + \psi(x))$

Using an antisymmetric activation function to learn faster

☆ Because The network is more balanced in handling positive and negative inputs to reducing bias, and helping to train data centered around an output of 0.

## IV. Normalizing the inputs



Mean Removal ⇒    Decorrelation ⇒    Covariance Equalization ⇒

1) Mean Removal ( Subtract the mean)
    i. Makes the input of the activation function in the range close to 0,
       where the gradient is larger.
         a. Accelerate the convergence rate of gradient descent
         b. Improve numerical stability to avoid gradient vanishing

    ii. Eliminate bias and let the NN focus on changes in the data.

2) Decorrelation ( PCA )
    i. Equalizes weight learning speed

    ii. Prevent overfitting by reduce the data dimensionality, decrease the number of parameter.
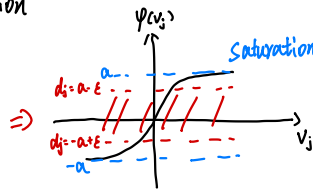
3) Covariance Equalization
    i. Normalize feature scales to ensure the weights are updated in a equalize magnitude,
       accelerating model convergence
    ii. Make the feature space more uniform.

☆ V. Weight initialization

We want $v_j$ falls within
the red range to achieve
a better gradient

$\Rightarrow$



$d_i = a - \varepsilon$

$d_j = -a + \varepsilon$

Saturation

$\varphi(v_j)$

$v_j$

$v_j = \sum_{i=1}^{m} w_{ji} \, y_i$

Considering $v_j$, $y_i$, $w_{ji}$ as random variables, we have:

$$\mu_v = \mathbb{E}[v_j] = \mathbb{E}\left[\sum_{i=1}^{m} w_{ji} \, y_i\right]$$

Assumption 1:

$\mu_y = \mathbb{E}[y_i] = 0$

$\sigma_y^2 = \mathbb{E}[(y_i - \mu_y)^2]$
$\quad\; = \mathbb{E}[y_i^2] = 1$

$$= \sum_{i=1}^{m} \mathbb{E}[w_{ji}] \, \mathbb{E}[y_i]$$

$$= 0$$

$$\sigma_v^2 = \mathbb{E}[(v_j - \mu_v)^2]$$
$$= \mathbb{E}[v_j^2]$$
$$= \mathbb{E}\left[\sum_{i=1}^{m} \sum_{k=1}^{m} w_{ji} \, w_{jk} \, y_i \, y_k\right]$$

$$= \sum_{i=1}^{m} \sum_{k=1}^{m} \mathbb{E}[w_{ji} w_{jk}] \underline{\mathbb{E}[y_i y_k]}$$

Assumption 2: Decorrelation step

$$\mathbb{E}[y_i y_k] = \begin{cases} 1 & , k = i \\ 0 & , k \neq i \end{cases}$$

$$= \sum_{i=1}^{m} \mathbb{E}[w_{ji}^2]$$

$$= m \sigma_w^2$$

if $\sigma_v = 1$. we have $\sigma_w^2 = \frac{1}{m}$

thus we randomly initialize weights with $\mu_w = 0$ and $\sigma_w^2 = \frac{1}{m}$.

reciprocal of the number of weight connection

③ Training workflow



Initialization $\Leftarrow w = \begin{cases} \mu_w = 0 \\ \sigma_w^2 = \frac{1}{m} \end{cases}$ ( Guaranteed gradient )

Input training set (Shuffle the epoch at each time)

Forward Computation

Backward Computation $① \delta_{j(n)}^{(L)} \begin{cases} e_j^{(L)}(n) \, \varphi_j'(v_j^{(L)}(n)) & \text{output } j \\ \\ \varphi_j(v_j^{(L)}(n)) \sum_k \delta_k^{(L+1)}(n) \, w_{kj}^{(L+1)}(n) & \text{hidden } j \end{cases}$

Converge?

② Adjust weights :

$$w_{ji}^{(L)}(n+1) = w_{ji}^{(L)}(n) + \alpha \Delta w_{ji}^{(L)}(n-1) + \eta \delta_j^{(L)}(n) y_i^{(L)}(n)$$

Done