hw 5

Yuanteng Chen

3039725444      1. Depthwise Separable Convolutions

(a). learnable parameters: $(3 \times 3 \times 3) \times 4 = 108$

(b)

$\begin{cases} \text{Depthwise} \\ \text{convolution}: \quad 3 \times 3 \times 3 = 27 \\ \\ \text{Pointwise} \\ \text{convolution}: \quad (1 \times 1 \times 3) \times 4 = 12 \end{cases}$

learnable parameters: $27 + 12 = 39$

2. Regularization and dropout

$$L(w) = ||y - Xw||_2^2 \quad (1)$$

$$L(\tilde{w}) = E_{R \sim Bernoulli(p)} \left[ ||y - (R \odot X)\tilde{w}||_2^2 \right] \quad (2)$$

$$L(w) = ||y - Xw||_2^2 + ||\Gamma w||_2^2 \quad (3)$$

(a) manipulate (2) to eliminate the expectations
and get, $L(\tilde{w}) = ||y - pX\tilde{w}||_2^2 + p(1-p)||\tilde{\Gamma}\tilde{w}||_2^2$

Solution,

assume $R \odot X = P$

$\therefore ||y - (R \odot X)\tilde{w}||_2^2 = ||y - P\tilde{w}||_2^2$

$= y^T y + \tilde{w}^T P^T P \tilde{w} - 2\tilde{w}^T P^T y$

$$\therefore E_{R \sim Bernoulli(p)} \left[ \| y - (R \odot X) \tilde{w} \|_2^2 \right]$$

$$= E_{R \sim B(p)} \left[ y^T y + w^T P^T P w - 2 w^T P^T y \right]$$

① $E_R [P]_{ij} = E_R [(R \odot X)_{ij}] = E_R [R_{ij}] \cdot X_{ij} = p X_{ij}$

② $E_R [2 w^T P^T y] = 2 p w^T x^T y$

③ $(E_R [(P^T P)])_{ij} = \sum_{k=1}^{N} E_R [R_{ki} R_{kj} X_{ki} X_{kj}]$

$$E_R [(P^T P)]_{ij} = \begin{cases} \sum_{k=1}^{N} E_R (R_{ki}) E_R (R_{kj}) \cdot X_{ki} \cdot X_{kj} = p^2 (x^T x)_{ij} & \\ & (i \neq j) \\ \\ \sum_{k=1}^{N} E_R [R_{ki}^2 X_{ki} X_{kj}] = \sum_{k=1}^{N} E_R [R_{ki}^2] X_{ki} X_{kj} & \\ \qquad\qquad = p (x^T x)_{ij} \quad (i = j) \end{cases}$$

$$(E_R [(P^T P)])_{ij} - p^2 (x^T x)_{ij} = \begin{cases} 0 & i \neq j \\ (p^2 - p)(x^T x)_{ij} & i = j \end{cases}$$

$$\therefore E_{R \sim B(p)} \left[ y^T y + w^T P^T P w - 2 w^T P^T y \right]$$

$$= y^T y - 2 p w^T x^T y + (p^2 w^T x^T x w - p^2 w^T x^T x w) + w^T E_R [P^T P] w$$

$$= (y^T y - 2 p w^T x^T y + p^2 w^T x^T x w) - p^2 w^T x^T x w + w^T E_R [P^T P] w$$

$$= \| y - p X w \|_2^2 + w^T (E_R [P^T P] - p^2 x^T x) w$$

$$= \| y - p X w \|_2^2 + w^T (p^2 - p) \, diag(x^T x) \, w$$

$\left( \text{only when } i = j, \ (E_R [(P^T P)])_{ij} - p^2 (x^T x)_{ij} = (p^2 - p)(x^T x)_{ij} \right)$

$$= \| y - p X w \|_2^2 + p(1-p) \| \tilde{\Gamma} w \|_2^2 \qquad \tilde{\Gamma} = \sqrt{diag(x^T x)}$$

(b). $L(\tilde{w}) = ||y - px\,\tilde{w}||_2^2 + p(1-p)||\tilde{\Gamma}\tilde{w}||_2^2$

assume $w = p\,\tilde{w}$

$L(\tilde{w}) = ||y - Xw||_2^2 + p(1-p)||\tilde{\Gamma}\frac{w}{p}||_2^2$

$\qquad = ||y - Xw||_2^2 + ||\sqrt{\frac{1-p}{p}}\,\tilde{\Gamma}w||_2^2$

$\qquad = ||y - Xw||_2^2 + ||\Gamma w||_2^2 \qquad (\Gamma = \sqrt{\frac{1-p}{p}}\,\tilde{\Gamma})$

(c) $\quad L(w) = ||y - Xw||_2^2 + ||\Gamma w||_2^2$

$\qquad\qquad \downarrow$

$\quad L(\tilde{w}) = ||y - \tilde{X}\tilde{w}||_2^2 + \lambda||\tilde{w}||_2^2$

Sol: assume $\tilde{w} = \Gamma w \quad \therefore w = \Gamma^{-1}\tilde{w}$

$\qquad \therefore \tilde{X}\Gamma = X \quad \tilde{X} = X \cdot \Gamma^{-1}$

3. Multiplicative Regularization beyond Dropout
expected training loss:
$\quad L(w) = E_{Rij \sim N(\mu, \sigma^2)}\left[||y - (R \odot X)w||_2^2\right]$

can be put in the form:
$\quad L(w) = ||y - \_(A)\_Xw||_2^2 + \_(B)\_||\Gamma w||_2^2$
where $\Gamma = (diag(X^TX))^{\frac{1}{2}}$

Sol:
in 2(a).

$E_{R \sim B(p)}\left[||y - (R \odot X)w||_2^2\right]$
$= ||y - pXw||_2^2 + p(1-p)||\Gamma w||_2^2 \qquad \Gamma = (diag(X^TX))^{\frac{1}{2}}$

in Bernoulli distribution:

$$P(X=k) = p^k (1-p)^{1-k}$$

$$E(X) = p, \qquad Var(X) = p(1-p)$$

in normal distribution

$$E(X) = \mu, \qquad Var(X) = \sigma^2$$

A: $\mu$      B: $\sigma^2$

## 4. Analyzing Distributed Training

| | Number of Message Sent | Size of each message |
|---|---|---|
| All-to-all | $n(n-1)$ | $p$ |
| Parameter Server | $2n$ | $p$ |
| Ring All-Reduce | $n(2(n-1))$ | $\dfrac{p}{n}$ |