

# CSC413 Homework 3

Name: Yulin Wang

Student #: 1003942326

## 1.1.1 Batch size vs. learning rate

As batch size increases, the gradient noise will decrease.

followed by a larger optimal learning rate.

If  $B$  is small initially, then the increase in the optimal learning rate will be more prominent.

## 1.1.2 Training steps vs. batch size

(a) Point C has the most efficient batch size.

For the batch sizes smaller than C, increasing batch size will

reduce the training time without much change in total compute.

For the batch sizes greater than C, the reduction in training time is minimal,

and would cost more compute.

(b) Point A : Regime: noise dominated

Point B : Regime: curvature dominated

## 1.1.3 Batch size, Optimizer, Normalization, Learning Rate

(a) options I and IV

(b) options II+ and III-

## 1.2 Model size, dataset size and compute

(a) (i) Model A has more parameters.

Since Model A converges to a smaller total loss in Figure 3.

(ii) At "X", Model B has been training for more iterations.

Since Model B has fewer parameters then it needs to train

more iterations to achieve the same total loss as model A.

(b) I will choose model A.

Since model A could converge to a smaller total loss than model B,

when given the same total compute.

## 2.1 Bias-variance decomposition

### 2.1.1

$$\begin{aligned} R(\hat{w}) &= E[(w_*^T \bar{x} - \hat{w}^T \bar{x})^2] \\ &= E[(w_*^T \bar{x} - \hat{w}^T \bar{x})^T (w_*^T \bar{x} - \hat{w}^T \bar{x})] \\ &= E[\bar{x}^T (w_* - \hat{w})(w_*^T - \hat{w}^T) \bar{x}] \\ &= E[\text{tr}[\bar{x}^T (w_* - \hat{w})(w_*^T - \hat{w}^T) \bar{x}]] \\ &= \text{tr}[E(\bar{x}^T \bar{x}) E[(w_* - \hat{w})(w_*^T - \hat{w}^T)]] \\ &= \text{tr}[E(\bar{x}^T \bar{x}) E(w_* w_*^T - w_* \hat{w}^T - \hat{w} w_*^T + \hat{w} \hat{w}^T)] \\ &\quad \textcircled{1} \quad \textcircled{2} \quad \textcircled{3} \quad \textcircled{4} \quad \textcircled{5} \end{aligned}$$

Since  $\hat{w} = (X^T X)^{-1} X^T (X w_* + \varepsilon) = w_* + (X^T X)^{-1} X^T \varepsilon$  when  $n > d$ .

Then ①  $E(\bar{x}^T \bar{x}) = Id$ ; ②  $E(w_* w_*^T) = \frac{d}{n} Id$

$$\begin{aligned} \textcircled{3} E(w_* \hat{w}^T) &= E[w_* (w_* + (X^T X)^{-1} X^T \varepsilon)^T] \\ &= E[w_* w_*^T + w_* \varepsilon^T X (X^T X)^{-1}] ; E(\varepsilon) = E(\varepsilon^T) = 0 \\ &= \frac{d}{n} Id + 0 = \frac{d}{n} Id \end{aligned}$$

$$\textcircled{4} \quad E(\hat{w}\hat{w}^T) = E[(\alpha u_* + \hat{\epsilon})^T (\alpha u_* + \hat{\epsilon})] = \frac{d}{n} I_d$$

$$\begin{aligned} \textcircled{5} \quad E(\hat{w}\hat{w}^T) &= E[(\alpha u_* + \alpha^T X^T \Sigma)(\alpha u_* + \alpha^T X^T \Sigma)^T] \\ &= E[(\alpha u_* + \alpha^T X^T \Sigma)(W_*^T + \Sigma^T X \alpha^T X^T)] \\ &= E(W_* W_*^T) + E(\alpha \alpha^T X^T \Sigma W_*^T) + E(X^T \Sigma^T \Sigma W_*^T) + E[\alpha^T X^T \Sigma^T \Sigma^T X \alpha^T X^T] \\ &= \frac{d}{n} I_d + 0 + 0 + E[\alpha^T X^T \Sigma^T \Sigma^T X \alpha^T X^T] \end{aligned}$$

$$\begin{aligned} \text{Thus, } R(\hat{w}) &= \text{tr}[E(\hat{w}\hat{w}^T) - W_* W_*^T - \hat{w}\hat{w}^T - \hat{w}\hat{w}^T] \\ &= \text{tr}[I_d (\frac{d}{n} I_d - \frac{d}{n} I_d - \frac{d}{n} I_d + \frac{d}{n} I_d + E[\alpha^T X^T \Sigma^T \Sigma^T X \alpha^T X^T])] \\ &= 0 + \text{tr} E[\alpha^T X^T \Sigma^T \Sigma^T X \alpha^T X^T]; \quad E[\Sigma \Sigma^T] = \sigma^2 I_n \\ &= 0 + \sigma^2 \cdot \text{tr}[X^T X] \end{aligned}$$

## 2.2 Deriving the exact expressions

### 2.2.1 ① under-parameterized case $n > d$

$$E[R(\hat{w})] = \frac{d\sigma^2}{n-d-1}$$

### ② over-parameterized case $n < d$ ; by 2.1.2

$$\begin{aligned} \text{tr}[I_d - X^T X] &= d-n; \quad \text{tr}[X^T X] = \frac{n}{d-n-1} \\ \Rightarrow E[R(\hat{w})] &= E[\frac{d}{n} \text{tr}[I_d - X^T X] + \sigma^2 \text{tr}[X^T X]] \\ \Rightarrow E[R(\hat{w})] &= \frac{d-n}{d} + \frac{n\sigma^2}{d-n-1} \end{aligned}$$

### 2.2.2 Double descent

#### (1) ① under-parameterized case $n > d$

$$\text{Set } E[R(\hat{w})] = 0 \Rightarrow \frac{d\sigma^2}{n-d-1} = 0 \Rightarrow \sigma = 0 \Rightarrow \text{variance equals to 0.}$$

#### ② over-parameterized case $n < d$

$$\text{Set } E[R(\hat{w})] = 0 \Rightarrow \frac{d-n}{d} + \frac{n\sigma^2}{d-n-1} = 0 \Rightarrow n=d \text{ and } (n=0 \text{ or } \sigma=0)$$

since  $n < d \Rightarrow$  contradicts  $\Rightarrow$  impossible

(2) NO, it does NOT always help generalization.

As the above formulas show, when adding more training examples might increase the variance.

There exists a trade-off between them.

## 2.3 Ridge regularization

2.3.2 ① When  $n$  increases, the model would be easier to have overfitting issue.

should increase  $\lambda$ .

② When  $\sigma$  increases, the model becomes less likely to have overfitting issue.

should decrease  $\lambda$ .

### 2.3.4

① From 2.2.1 under:  $E[R(\hat{w})] = \frac{d\sigma^2}{n-d-1}$

$$\text{over: } E[R(\hat{w})] = \frac{d-n}{d} + \frac{n\sigma^2}{d-n-1}$$

$$\text{ridge-regularized: } E[R(\hat{w}_{\text{reg}})] = \sigma^2 \cdot \frac{-(1-\gamma+\lambda) + \sqrt{(1-\gamma+\lambda)^2 + 4\gamma\lambda}}{2\lambda}$$

Instead of separating  $d$  and  $n$ , the ridge-regularized estimator generalizes than by using  $\gamma = \frac{d}{n}$ .

From the formula, as  $\gamma = \frac{d}{n}$  increases,  $E[R(\hat{w}_{\text{reg}})]$  increases, which makes it easier to find effects by conditions.

② YES, adding more training data always leads to better test performance, given  $\lambda = \sigma^2\gamma$ .