

1.1.1) As batch size increases, optimal learning rate should increase too because larger batch size leads to more confident gradients (more confident in the direction of the descent) (i.e. less stochastic) (the minibatch gradient noise is proportional to  $\frac{1}{B}$  of stochastic noise). so we can use larger learning rate,

$$\begin{aligned} \text{Var}[g_B(w)] &= \text{Var}\left[\frac{1}{B} \sum_{i=1}^B g_i(w)\right] = \frac{1}{B^2} \sum_{i=1}^B \text{Var}[g_i(w)] = \frac{1}{B^2} \sum_{i=1}^B \text{Var}[\nabla L(w) + \varepsilon_i] \\ &= \frac{1}{B^2} \sum_{i=1}^B \text{Var}[\varepsilon_i] = \frac{1}{B^2} B \sigma^2 \\ &= \frac{1}{B} \sigma^2 \end{aligned}$$

1.1.2) a) point C has the most efficient batch size because for batch size smaller than C, increasing batch size significantly decreases the number of steps needed to reach the validation error. For batch size larger than C, this convergence benefits see diminishing returns, but the time per weight update should increase with more batch size.

b) Point A: Regime noise dominated

Point B: Regime curvature dominated.

1.1.3)

a) I, IV

b) II+, III-

1.2

a)

i) Model A has more number of parameters than model B because as total compute (number of training steps) increase, the test loss for B is able to reach much lower than test loss for A, meaning the increased complexity of B makes it able to generate better fit to test data.

ii) at X, total compute A = total compute B, and since the number of parameters of A > number of parameters B,  
⇒ the number of training steps must be more for B than A to use the same amount of total compute.

b)

If I have an urgent deadline coming up, I would use the larger model (A) since it requires less training steps (less wall clock time) to reach the same test loss.  
On the other hand, if I do not have access to good machines, it might be better to train a smaller model assuming it can reach the level of test loss wanted.

2.1.1)

$$R(\hat{W}) = E_{\mathbf{x}, \mathbf{\Sigma}, \mathbf{W}_*} [(\mathbf{W}_*^\top \tilde{\mathbf{x}} - \hat{\mathbf{W}}^\top \tilde{\mathbf{x}})^2]$$

since  $\tilde{\mathbf{x}}, \mathbf{\Sigma}, \mathbf{N}_*$  are independent we are able to pull variables out of the expectations over other variables.

$$= E_{\mathbf{x}, \mathbf{\Sigma}, \mathbf{W}_*} [(\mathbf{W}_*^\top \tilde{\mathbf{x}})^2 - 2 \mathbf{W}_*^\top \tilde{\mathbf{x}} \hat{\mathbf{W}}^\top \tilde{\mathbf{x}} + (\hat{\mathbf{W}}^\top \tilde{\mathbf{x}})^2]$$

$$= \underbrace{E_{\mathbf{x}, \mathbf{\Sigma}, \mathbf{W}_*} [(\mathbf{W}_*^\top \tilde{\mathbf{x}})^2]}_{\textcircled{1}} - 2 \underbrace{E_{\mathbf{x}, \mathbf{\Sigma}, \mathbf{W}_*} [\mathbf{W}_*^\top \tilde{\mathbf{x}} \hat{\mathbf{W}}^\top \tilde{\mathbf{x}}]}_{\textcircled{2}} + \underbrace{E_{\mathbf{x}, \mathbf{\Sigma}, \mathbf{W}_*} [(\hat{\mathbf{W}}^\top \tilde{\mathbf{x}})^2]}_{\textcircled{3}}$$

by linearity of expectation.

next we evaluate each of the 3 terms separately:

$$\textcircled{1} E_{\mathbf{x}, \mathbf{\Sigma}, \mathbf{W}_*} [(\mathbf{W}_*^\top \tilde{\mathbf{x}})^2] = E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [E_{\mathbf{x}} [(\mathbf{W}_*^\top \tilde{\mathbf{x}})^2]] \right]$$

$$= E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} \left[ \left( E_{\mathbf{x}} [\mathbf{W}_*^\top \tilde{\mathbf{x}}] \right)^2 + \text{Var}_{\mathbf{x}} [\mathbf{W}_*^\top \tilde{\mathbf{x}}] \right] \right]$$

$$= E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} \left[ \mathbf{0}^2 + \mathbf{W}_*^\top \text{Var}_{\mathbf{x}} [\tilde{\mathbf{x}}] \mathbf{W}_* \right] \right]$$

$$= E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [\mathbf{W}_*^\top \mathbf{W}_*] \right]$$

$$= E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [\text{Tr}(\mathbf{W}_*^\top \mathbf{W}_*)] \right]$$

$$= E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [\text{Tr}(\mathbf{W}_* \mathbf{W}_*)] \right] \quad \text{by cyclic property of trace.}$$

$$= E_{\mathbf{\Sigma}} \left[ \text{Tr}(E_{\mathbf{W}_*} [\mathbf{W}_* \mathbf{W}_*^\top]) \right]$$

$$= E_{\mathbf{\Sigma}} \left[ \text{Tr} \left( \frac{1}{d} \mathbf{I}_d \right) \right]$$

$$= E_{\mathbf{\Sigma}} [1]$$

$$= 1$$

$$\textcircled{2} -2 E_{\mathbf{x}, \mathbf{\Sigma}, \mathbf{W}_*} [\mathbf{W}_*^\top \tilde{\mathbf{x}} \hat{\mathbf{W}}^\top \tilde{\mathbf{x}}] = -2 E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [E_{\mathbf{x}} [\mathbf{W}_*^\top \tilde{\mathbf{x}} \hat{\mathbf{W}}^\top \tilde{\mathbf{x}}]] \right]$$

$$= -2 E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [E_{\mathbf{x}} [\mathbf{W}_*^\top \tilde{\mathbf{x}} \mathbf{t}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \tilde{\mathbf{x}}]] \right] \quad \text{since } \hat{\mathbf{W}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} \Rightarrow \hat{\mathbf{W}}^\top = \mathbf{t}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}$$

$$= -2 E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [E_{\mathbf{x}} [\mathbf{W}_*^\top \tilde{\mathbf{x}} (\mathbf{w}_*^\top \mathbf{X}^\top \mathbf{\Sigma}^\top) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \tilde{\mathbf{x}}]] \right] \quad \text{since } \mathbf{t} = \mathbf{X} \mathbf{w}_* + \mathbf{\Sigma} \Rightarrow \mathbf{t}^\top = \mathbf{w}_*^\top \mathbf{X}^\top \mathbf{\Sigma}^\top$$

$$= -2 E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [E_{\mathbf{x}} [\mathbf{W}_*^\top \tilde{\mathbf{x}} \mathbf{w}_*^\top \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \tilde{\mathbf{x}}]] + E_{\mathbf{\Sigma}} [E_{\mathbf{W}_*} [\mathbf{W}_*^\top \tilde{\mathbf{x}} \mathbf{\Sigma}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \tilde{\mathbf{x}}]] \right]$$

$$= -2 E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [E_{\mathbf{x}} [(\mathbf{W}_*^\top \tilde{\mathbf{x}})^2]] \right] - 2 E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [E_{\mathbf{x}} [\text{tr}(\mathbf{W}_*^\top \tilde{\mathbf{x}} \mathbf{\Sigma}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \tilde{\mathbf{x}})]] \right]$$

$$= -2(1) - 2 E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [E_{\mathbf{x}} [\text{tr}(\mathbf{W}_*^\top \tilde{\mathbf{x}} \mathbf{\Sigma}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \tilde{\mathbf{x}})]] \right] \quad \text{since } E_{\mathbf{\Sigma}} [E_{\mathbf{W}_*} [E_{\mathbf{x}} [(\mathbf{W}_*^\top \tilde{\mathbf{x}})^2]]] = 1 \text{ by } \textcircled{1}$$

$$= -2 - 2 E_{\mathbf{\Sigma}} \left[ E_{\mathbf{W}_*} [E_{\mathbf{x}} [\text{tr}(\mathbf{\Sigma}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \tilde{\mathbf{x}} \mathbf{W}_*^\top \tilde{\mathbf{x}})]] \right] \quad \text{by cyclic property of trace}$$

$$\begin{aligned}
&= -2 - 2 \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \text{tr} (\Sigma X(XX^{-1}) \tilde{x} W_*^\top \tilde{x}) \right] \right] \\
&= -2 - 2 \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \text{tr} (\Sigma X(XX^{-1}) \mathbb{E}_{\tilde{x}} [\tilde{x} W_*^\top \tilde{x}]) \right] \right] \\
&= -2 - 2 \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \text{tr} (\Sigma X(XX^{-1}) I_d W_*) \right] \right] \quad \text{since } \mathbb{E}_{\tilde{x}} [\tilde{x} W_*^\top \tilde{x}] = I_d W_* \\
&= -2 - 2 \mathbb{E}_{\Sigma} \left[ \text{tr} (\Sigma X(XX^{-1}) I_d \mathbb{E}_{W_*} [W_*]) \right] \\
&= -2 - 2 \text{tr} (\mathbb{E}_{\Sigma} [\Sigma] X(XX^{-1}) I_d \mathbb{E}_{W_*} [W_*]) \\
&= -2 - 0 \\
&= -2
\end{aligned}$$

③  $\mathbb{E}_{\tilde{x}, \Sigma, W_*} [(\hat{w}^\top \tilde{x})^2] = \mathbb{E}_{\Sigma W_* \tilde{x}} \left[ \mathbb{E}_{\tilde{x}} [(\hat{w}^\top X(XX^{-1}) \tilde{x})^2] \right]$  since  $\hat{w}^\top = t^\top X(XX^{-1})$

$$\begin{aligned}
&= \mathbb{E}_{\Sigma W_* \tilde{x}} \left[ \mathbb{E}_{\tilde{x}} \left[ (\mathbb{E}_{\tilde{x}} [t^\top X(XX^{-1}) \tilde{x}])^2 + \text{Var}[t^\top X(XX^{-1}) \tilde{x}] \right] \right] \\
&= \mathbb{E}_{\Sigma W_* \tilde{x}} \left[ \mathbb{E}_{\tilde{x}} \left[ (\mathbb{E}_{\tilde{x}} [(\omega_*^\top X^\top + \varepsilon) X(XX^{-1}) \tilde{x}])^2 + \text{Var}[(\omega_*^\top X^\top + \varepsilon) X(XX^{-1}) \tilde{x}] \right] \right] \quad \text{since } t^\top = \omega_*^\top X^\top + \varepsilon \\
&\quad \cancel{\text{0}} \\
&= \mathbb{E}_{\Sigma W_* \tilde{x}} \left[ \mathbb{E}_{\tilde{x}} \left[ ((\omega_*^\top X^\top + \varepsilon) X(XX^{-1}) \mathbb{E}_{\tilde{x}} [\tilde{x}])^2 + \text{Var}[(\omega_*^\top X^\top + \varepsilon) X(XX^{-1}) \tilde{x}] \right] \right] \\
&= \mathbb{E}_{\Sigma W_* \tilde{x}} \left[ \mathbb{E}_{\tilde{x}} \left[ 0 + \text{Var}[(\omega_*^\top X^\top + \varepsilon) X(XX^{-1}) \tilde{x}] \right] \right] \\
&= \mathbb{E}_{\Sigma W_* \tilde{x}} \left[ \mathbb{E}_{\tilde{x}} \left[ \text{Var}[\omega_*^\top X^\top X(XX^{-1}) \tilde{x} + \varepsilon X(XX^{-1}) \tilde{x}] \right] \right] \\
&= \mathbb{E}_{\Sigma W_* \tilde{x}} \left[ \mathbb{E}_{\tilde{x}} \left[ \text{Var}[\omega_*^\top \tilde{x} + \varepsilon X(XX^{-1}) \tilde{x}] \right] \right] \\
&= \mathbb{E}_{\Sigma W_* \tilde{x}} \left[ \mathbb{E}_{\tilde{x}} \left[ \text{Var}[(\omega_*^\top + \varepsilon X(XX^{-1})) \tilde{x}] \right] \right] \\
&= \mathbb{E}_{\Sigma W_* \tilde{x}} \left[ \mathbb{E}_{\tilde{x}} \left[ (\omega_*^\top + \varepsilon X(XX^{-1})) \text{Var}[\tilde{x}] (\omega_*^\top + \varepsilon X(XX^{-1}))^\top \right] \right] \\
&= \mathbb{E}_{\Sigma W_* \tilde{x}} \left[ \mathbb{E}_{\tilde{x}} \left[ (\omega_*^\top + \varepsilon X(XX^{-1})) (\omega_*^\top + (XX^{-1})^\top \varepsilon) \right] \right] \\
&= \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \omega_*^\top \omega_* + \omega_*^\top (XX^{-1})^\top \varepsilon + \varepsilon X(XX^{-1}) \omega_* + \varepsilon X(XX^{-1}) (XX^{-1})^\top \varepsilon \right] \right] \\
&\quad \cancel{\text{0}} \\
&= \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \omega_*^\top \omega_* \right] + \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \omega_*^\top (XX^{-1})^\top \varepsilon \right] \right] + \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \varepsilon X(XX^{-1}) \omega_* \right] \right] + \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \varepsilon X(XX^{-1}) (XX^{-1})^\top \varepsilon \right] \right] \right] \\
&= \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \omega_*^\top \omega_* \right] + \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \omega_*^\top (XX^{-1})^\top \varepsilon \right] \right] + \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \varepsilon X(XX^{-1}) \omega_* \right] \right] + \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \varepsilon X(XX^{-1}) (XX^{-1})^\top \varepsilon \right] \right] \right] \\
&= \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \omega_*^\top (XX^{-1})^\top \varepsilon \right] \right] + \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \varepsilon X(XX^{-1}) \omega_* \right] \right] + \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \varepsilon X(XX^{-1}) (XX^{-1})^\top \varepsilon \right] \right] \\
&= \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \omega_*^\top (XX^{-1})^\top \varepsilon \right] \right] + \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \varepsilon X(XX^{-1}) \omega_* \right] \right] + \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \varepsilon X(XX^{-1}) (XX^{-1})^\top \varepsilon \right] \right] \\
&= \mathbb{E}_{\Sigma W_*} \left[ \mathbb{E}_{\tilde{x}} \left[ \varepsilon X(XX^{-1}) (XX^{-1})^\top \varepsilon \right] \right]
\end{aligned}$$

by second half of ①

$$\begin{aligned}
&= 1 + \mathbb{E}_{\boldsymbol{\varepsilon}} [\text{tr}(\boldsymbol{\varepsilon} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}^T \boldsymbol{\varepsilon})] \\
&= 1 + \mathbb{E}_{\boldsymbol{\varepsilon}} [\text{tr}(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}^T))] \quad \text{by cyclic property of trace} \\
&= 1 + \text{tr} \left( \mathbb{E}_{\boldsymbol{\varepsilon}} [\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}^T)] \right) \\
&= 1 + \text{tr} \left( \mathbb{E}_{\boldsymbol{\varepsilon}} [\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}^T) \right) \\
&= 1 + \sigma^2 \text{tr} (\mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}^T)) \\
&= 1 + \sigma^2 \text{tr} (\mathbf{X}^T \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}) \quad \text{by cyclic property of trace} \\
&= 1 + \sigma^2 \text{tr} (\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}) \\
&= 1 + \sigma^2 \text{tr} ((\mathbf{X}^T \mathbf{X})^{-1})
\end{aligned}$$

$$\begin{aligned}
\text{so } R(\hat{\mathbf{w}}) &= \mathbb{E}_{\mathbf{X}, \boldsymbol{\varepsilon}, \mathbf{W}_*} [(\mathbf{W}_*^T \tilde{\mathbf{z}} - \hat{\mathbf{w}}^T \tilde{\mathbf{z}})^2] = \textcircled{1} + \textcircled{2} + \textcircled{3} = 1 - 2 + 1 + \sigma^2 \text{tr} ((\mathbf{X}^T \mathbf{X})^{-1}) \\
&= 0 + \sigma^2 \text{tr} ((\mathbf{X}^T \mathbf{X})^{-1})
\end{aligned}$$

as wanted //

2.2.1)

$n > d+1$ :

$$\begin{aligned} E[R(\hat{\mathbf{w}})] &= E\left[\sigma^2 \text{Tr}((\mathbf{X}^\top \mathbf{X})^{-1})\right] \\ &= \sigma^2 E\left[\text{Tr}((\mathbf{X}^\top \mathbf{X})^{-1})\right] \\ &= \sigma^2 \frac{d}{n-d-1} \end{aligned}$$

$n < d+1$ :

$$\begin{aligned} E[R(\hat{\mathbf{w}})] &= E\left[\frac{1}{d} \text{Tr}(\mathbf{I}_d - \mathbf{X}(\mathbf{X}^\top)^{-1}\mathbf{X}) + \sigma^2 \text{Tr}((\mathbf{X}\mathbf{X}^\top)^{-1})\right] \\ &= \frac{1}{d} E\left[\text{Tr}(\mathbf{I}_d - \mathbf{X}(\mathbf{X}^\top)^{-1}\mathbf{X})\right] + \sigma^2 E\left[\text{Tr}((\mathbf{X}\mathbf{X}^\top)^{-1})\right] \\ &= \frac{1}{d} E[d-n] + \sigma^2 E\left[\text{Tr}((\mathbf{G}^\top \mathbf{G})^{-1})\right] \\ &= \frac{1}{d}(d-n) + \sigma^2 \frac{n}{d-n-1} \\ &= 1 - \frac{n}{d} + \sigma^2 \frac{n}{d-n-1} \end{aligned}$$

take  $\mathbf{G} = \mathbf{X}^\top \in \mathbb{R}^{d \times n}$   
 $\Rightarrow \mathbf{G}^\top = \mathbf{X} \in \mathbb{R}^{n \times d}$

2.2.2)

1)

if  $n < d+1$ , then  $\frac{n}{d} < 1$   
 $\Rightarrow 1 - \frac{n}{d} + \sigma^2 \frac{n}{d-n-1} > 0$ , since  $1 - \frac{n}{d} > 0$ . This means  $E[R(\hat{\mathbf{w}})] > 0$  for all  $n < d+1$ .

if  $n > d+1$ , then  $E[R(\hat{\mathbf{w}})] = \sigma^2 \frac{d}{n-d-1} = 0 \Rightarrow$  ①  $n = \infty$ ,  $d < \infty$ ,  $\sigma < \infty$

OR ②  $\sigma = 0$

OR ③  $d = 0$ .

so  $E[R(\hat{\mathbf{w}})] = 0$  only if  $n > d+1$  and

①  $n = \infty$ ,  $d < \infty$ ,  $\sigma < \infty \Rightarrow$  inf training data.

OR ②  $\sigma = 0 \Rightarrow$  no noise

OR ③  $d = 0 \Rightarrow$  meaningless in context

2) No, adding more training examples does not always help generalization.  
For example, in the overparametrized case,  $\frac{\partial E[R(\hat{\mathbf{w}})]}{\partial n} = -\frac{1}{d} + \frac{\sigma^2 n}{(d-n-1)^2} + \frac{\sigma^2}{d-n-1}$ ,

and this can be positive for certain  $\sigma, n, d$  (for example, take  $\sigma=1, n=1, d=4$ , then  $E[R(\hat{\mathbf{w}})] = 1 - \frac{1}{4} + 1^2 \frac{1}{4-1} = 1.25$ ; but same  $\sigma, d$  with increased  $n=2$  gives  $E[R(\hat{\mathbf{w}})] = 1 - \frac{2}{4} + 1^2 \frac{2}{4-2-1} = 2.5$ , which is larger.)

2.3.2)

The larger the sample size  $n$ , the less prone the model is to overfitting, so we should take smaller  $\lambda$ . On the other hand, larger noise in the dataset could cause the model to overfit on the noise, so we should take larger  $\lambda$  to regularize the model.

Based on Figure 4,

2.3.4) The test loss of the unregularized estimator  $E[R(\hat{w})]$  in 2.2.1 decreases at first as sample size increases until around  $n=400$  and then increases suddenly until  $n=500$  before dropping down. On the other hand the ridge-regularized  $E[R(\hat{w}_\lambda)]$  test loss decreases stably as sample size  $n$  increases. (The two have similar test loss for  $n < 400$  and  $n > 600$ .)

Based on Figure 4, under appropriate ridge regularization ( $\lambda = \sigma^2 \gamma$ ), adding more training data always leads to better test performance.