

Part a

$$\begin{aligned}
 \text{i.i} \quad \text{MSE}(\hat{\omega}) &= E_T \left[\| \hat{\omega} - \omega^* \|_2^2 \right] \\
 &\geq E_T \left[(\hat{\omega} - \omega)^T (\hat{\omega} - \omega^*) \right] \\
 &\geq E_T \left\{ \left[\hat{\omega} - \mu_{\hat{\omega}} + \mu_{\hat{\omega}} - \omega^* \right]^T \left(\hat{\omega} - \mu_{\hat{\omega}} + \mu_{\hat{\omega}} - \omega^* \right) \right\} \\
 &\geq E_T \left[(\hat{\omega} - \mu_{\hat{\omega}})^T (\hat{\omega} - \mu_{\hat{\omega}}) + 2(\hat{\omega} - \mu_{\hat{\omega}})^T (\mu_{\hat{\omega}} - \omega^*) \right. \\
 &\quad \left. + (\mu_{\hat{\omega}} - \omega^*)^T (\mu_{\hat{\omega}} - \omega^*) \right] \\
 &\geq E_T \left[\| \hat{\omega} - \mu_{\hat{\omega}} \|_2^2 \right] + 2E_T \left[(\hat{\omega} - \mu_{\hat{\omega}})^T (\mu_{\hat{\omega}} - \omega^*) \right. \\
 &\quad \left. + \| \mu_{\hat{\omega}} - \omega^* \|_2^2 \right] \\
 &\geq E_T \left[\| \hat{\omega} - \mu_{\hat{\omega}} \|_2^2 \right] + 2(\mu_{\hat{\omega}} - \omega^*)^T E_T (\hat{\omega} - \mu_{\hat{\omega}}) \\
 &\quad + \| \mu_{\hat{\omega}} - \omega^* \|_2^2 \\
 &\geq E_T \left[\| \hat{\omega} - \mu_{\hat{\omega}} \|_2^2 \right] + 0 + \| \mu_{\hat{\omega}} - \omega^* \|_2^2 \\
 &\geq \text{Var}(\hat{\omega}) + \| \mu_{\hat{\omega}} - \omega^* \|_2^2 \\
 &\geq \text{tr}(\text{cov}(\hat{\omega})) + \| \mu_{\hat{\omega}} - \omega^* \|_2^2
 \end{aligned}$$

Part b.

1.2

$$\hat{w}_{OLS} = (X^T X)^{-1} X^T y$$

$$E_T[\hat{w}_{OLS}] = E[(X^T X)^{-1} X^T y]$$

$X, y \rightarrow$ Random Variables

Given model $y = Xw + \epsilon$

$\epsilon \rightarrow$ normal distribution

$$E[y] = Xw$$

So,

$$E[\hat{w}_{OLS}] = (X^T X)^{-1} X^T E[y]$$

$$= (X^T X)^{-1} X^T X w = w$$

hence,

OLS estimator \rightarrow unbiased of w

$$2. \quad \hat{w}_{Ridge} = (X^T X + \lambda I_p)^{-1} X^T y$$

$$E_T[\hat{w}_{Ridge}] = (X^T X + \lambda I_p)^{-1} X^T E[y]$$

$$= (X^T X + \lambda I_p)^{-1} X^T X w$$

$\Rightarrow \hat{w}_{Ridge}$ is biased of w

Part-c

1.3

$$\hat{w}_{\text{ridge}} = (x^T x + \sigma^2 I_p)^{-1} x^T y$$

$$\text{Cov}(\hat{w}_{\text{ridge}}) = \sigma^2 (x^T x + \sigma^2 I_p)^{-1}$$

trace of matrix \geq Sum of eigenvalues lets
 $\lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \geq \lambda_p$ are eigenvalues of matrix

$x^T x$. Then eigenvalues of $x^T x + \sigma^2 I_p$ are $\lambda_i + \sigma^2$
 for $i, 1, 2, \dots, p$

$$\text{So, } \text{tr}[\text{Cov}(\hat{w}_{\text{ridge}})] = \sigma^2 \sum_{i=1}^p \frac{1}{\lambda_i + \sigma^2}$$

Similarly, $\hat{w}_{OLS} = (x^T x)^{-1} x^T y$

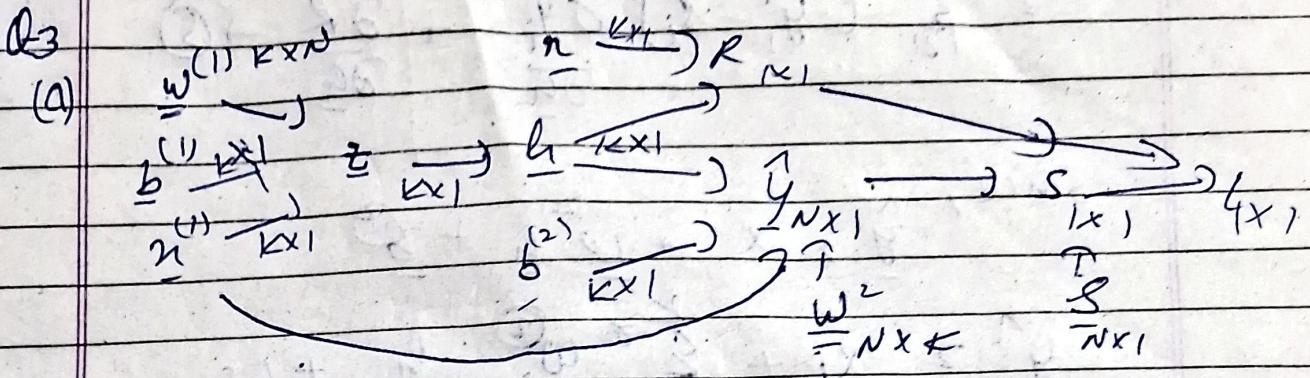
$$\text{Cov}(\hat{w}_{OLS}) = \sigma^2 (x^T x)^{-1}$$

trace of this Cov matrix \geq Sum of reciprocals of eigenvalues
 of $x^T x$

$$\text{tr}[\text{Cov}(\hat{w}_{OLS})] \geq \sigma^2 \sum_{i=1}^p \frac{1}{\lambda_i}$$

$$\lambda_i + \sigma^2 > \lambda_i \text{ so}$$

$$\text{tr}[\text{Cov}(\hat{w}_{\text{ridge}})] < \text{tr}[\text{Cov}(\hat{w}_{OLS})]$$



(b) $\bar{v} = \frac{\partial L}{\partial v} \rightarrow \text{rotation}$

→ Backward pass:

$$\underline{L} = \frac{\partial L}{\partial \underline{L}} = 1 \rightarrow \textcircled{1}$$

$$\bar{s} = \frac{\partial L}{\partial s} = 1 \rightarrow \textcircled{2}$$

$$\{ A_s, L = S + R \}$$

$$\hat{\underline{y}} = \frac{\partial L}{\partial \hat{\underline{y}}} = \frac{\partial L}{\partial s} \frac{\partial s}{\partial \hat{\underline{y}}} = 1 - \frac{\partial s}{\partial \hat{\underline{y}}} \rightarrow \textcircled{3}$$

$$S = \frac{1}{2} \| \hat{\underline{y}} - \underline{s} \|_2^2 = \frac{1}{2} (\hat{\underline{y}} - \underline{s})^T (\hat{\underline{y}} - \underline{s})$$

$$\frac{\partial s}{\partial \hat{\underline{y}}} = \hat{\underline{y}} - \underline{s}$$

$$\text{So, } \hat{\underline{y}} = \underline{s} - \underline{s} = \underline{0} \rightarrow \textcircled{4}$$

$$\underline{s} = \frac{\partial L}{\partial \underline{s}} = \sum_{i=1}^N \frac{\partial L}{\partial \hat{\underline{y}}_i} \frac{\partial \hat{\underline{y}}_i}{\partial \underline{s}} + \frac{\partial L}{\partial R} \frac{\partial R}{\partial \underline{s}}$$

$$\frac{\partial L}{\partial \hat{\underline{y}}_i} = \underline{g}_i - \underline{s}_i$$

$$= \sum_{i=1}^N (\hat{y}_i - \delta_i) \frac{\partial \hat{y}_i}{\partial h_i} + \frac{\partial L}{\partial h} \rightarrow \textcircled{S}$$

Since, $\hat{y} = x + w_0 h + b$

$$\hat{y}_i = x_i + \sum_{k=1}^K w_{ik} h_k + b_i$$

$$\frac{\partial \hat{y}_i}{\partial h_i} = \begin{bmatrix} w_{i0} \\ w_{i1} \\ \vdots \\ w_{iK} \end{bmatrix} \Rightarrow \sum_{i=1}^N (\hat{y}_i - \delta_i) \cdot \frac{\partial \hat{y}_i}{\partial h_i} \\ = \hat{w}^T (\hat{y} - \underline{\delta})$$

Since, $R_2 \sum_{i=1}^N \delta_i = \sum \epsilon_i \delta_i$

$$\Rightarrow \frac{\partial L}{\partial h} = \underline{\epsilon}$$

So, $\underline{\delta} = \hat{w}^T (\hat{y} - \underline{\delta}) + \underline{\epsilon} \rightarrow \textcircled{Q}$

$$\underline{\epsilon} = \frac{\partial L}{\partial \underline{z}} = \frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial \underline{z}}, \quad \{ \because h = g(z) \}$$

$$\therefore \underline{\epsilon} = \hat{w}^T (\hat{y} - \underline{\delta}) + \underline{\epsilon} : \sigma'(z) \rightarrow \textcircled{P}$$

$$\hat{u} = \frac{\partial L}{\partial \underline{u}} + \sum_{i=1}^k \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial \underline{u}} + \sum_{i=1}^n \frac{\partial L}{\partial e_i} \frac{\partial e_i}{\partial \underline{u}}$$

1st term :-

$$z = \underline{w}^T \underline{x} + b$$

$$z_i = \sum_{j=1}^N w_{ij} x_j + b_i$$

$$\frac{\partial z_i}{\partial \underline{u}} = \begin{bmatrix} w_{i1} \\ \vdots \\ w_{iN} \end{bmatrix} \Rightarrow \sum \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial \underline{u}}$$

$$= \underline{w}^T \frac{\partial L}{\partial \underline{z}}$$

2nd Term :-

$$\hat{y} = \underline{x} + \underline{w}^T \underline{h} + b$$

$$\hat{y}_i = x_i + \sum_{j=1}^k w_{ij} h_j + b_i$$

$$\frac{\partial \hat{y}_i}{\partial \underline{u}} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ i \\ 0 \end{bmatrix} \Rightarrow \sum_{i=1}^N \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial \underline{u}}$$

$$2 \sum_{i=1}^N (y_i - \hat{y}_i) \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$= \underline{y} - \underline{\hat{y}}$$

$$\underline{n} = \underline{w}^T \left[\left(\underline{w}^T (\underline{y} - \underline{\hat{y}}) + \varepsilon \right) \cdot \sigma'(\underline{z}) \right] + \underline{\hat{y}} - \underline{\delta}$$

(Q4.)

(i) Because of mini-batch gradient descent. The update of parameters happen based on only one batch ~~only~~ in each iteration. So, our data suffers a drop in accuracy.

(ii) 91.59%.

(iii) input layer $\rightarrow 784 + \frac{1}{(\text{bias})} = 785$

output layer $\rightarrow 10$

total parameter $= 784 \times 10$
 $= 7840$

(iv) Epoch \rightarrow one pass over training data
No. of iterations in each pass

\Rightarrow training data size $= \frac{60000}{100}$

$= 600$

\Rightarrow Loop runs 2000 times

so, total epochs $= \frac{2000}{600}$

$= 3.33$

hence,

No. of epochs $= 3$
We completed

(V)

Accuracy 94.78%

Yes, it's improved over single-layer neural network.

(Vi)

No. of hidden unit = 10

Test accuracy = 91.55%

→ No. of hidden unit = 400

Test accuracy = 94.78%

No. of hidden unit = 600

Test accuracy = 95.05% → With increase in hidden units in ^{4th} layer, increase in test accuracy.

(Vii)

Test accuracy drops to 75.62% .

(Plot)

(Viii)

For fix layers network:

Test Accuracy = 88.37% .which is ~~less than~~ 2-layer network with random property.

(Plot)

(IX) Test accuracy is 97.02% which is more than sigmoid.

Plot

(X) Test accuracy is 11.35%.

Plot

(XI) There is no NaN values within 2000 iterations. But more than 2000 iterations, NaN values start appearing in training set. Because of cross-entropy function, in which directly log of h_i is taken. So, when c_i is 0 it goes to undefined which might happen in ReLU activation function with negative logit values.

→ The F-cross-entropy function adds a similar epsilon before taking log. So, there is no effect on changing F-softmax & returning just the logit has any effect on this.

Therefore,

By changing cross-entropy function to F-cross-entropy function NaN values stopped to appearing.

Plot

(xii)

$$\begin{aligned}
 \text{Input layer} &\rightarrow 784 + 1 = 785 \\
 \text{hidden layer 1} &\rightarrow 200 + 1 = 201 \\
 2 &\rightarrow 100 + 1 = 101 \\
 3 &\rightarrow 60 + 1 = 61 \\
 4 &\rightarrow 30 + 1 = 31
 \end{aligned}$$

Output layer = 10

$$\begin{aligned}
 \text{total parameters} &= 785 \times 200 + 201 \times 100 + \\
 &101 \times 60 + 61 \times 30 + 31 \times 10 \\
 &= 785300
 \end{aligned}$$

(xiii) With optim Adam, learning rate = 0.03,
 ReLU activations, random initialization of weights
 & 0.1 initialization of biases.
 Test accuracy = 7.45%
 which is more than optim S.G.D. with ReLU
 activations, with 0.1 learning rate.

Plot

There are 2 changes
 we use lower learning rate with use of Adam.
 which adjust the size of jump we need to make.
 So, it will ~~slow down~~ if we are moving in
 some direction. This leads to higher accuracy.

(xiv) Test accuracy = 94.06%.

There is no improvement over 2000 iteration in
 accuracy; so, increase in more iterations will not
 improve accuracy.

Q.S.

(i) Hidden layer in 3rd layer = $12 \times 7 \times 7$
 $= 588$

↳ flat \rightarrow total no. of hidden units.

(ii) With 5 layer network & size of last 2 layers being 200 and 100.

Test accuracy = 98.81%
Plot

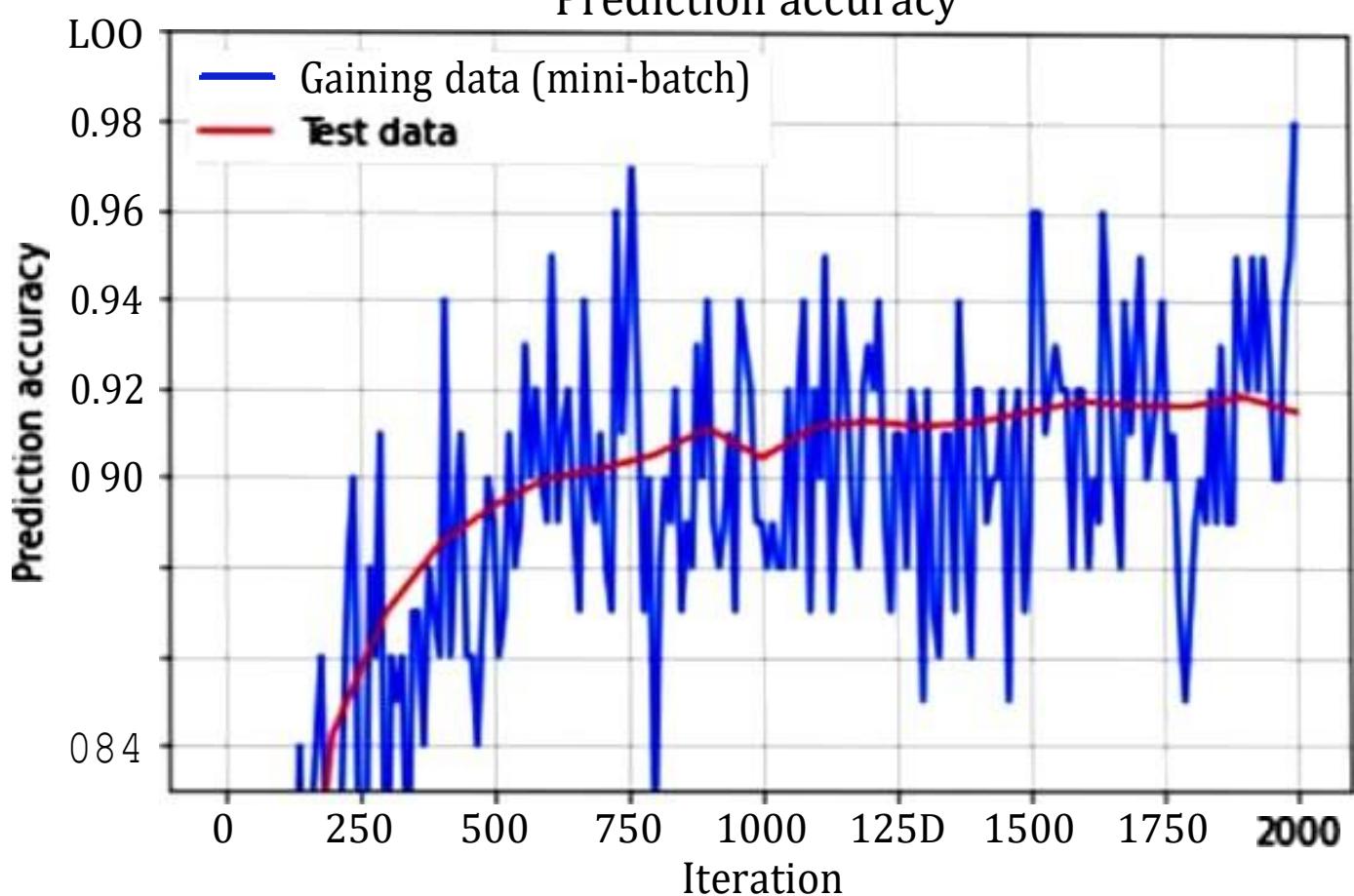
(iii) Test accuracy = 99.10%
Plot

(iv) At last, test accuracy is 100%. & saturated. Hence, now model is starting to overfit the data. The training accuracy was not tending to 100%. In earlier neural network w/o CNN. In CNN due to increased in complexity of convolutional & dense layers, the neural network has started overfitting.

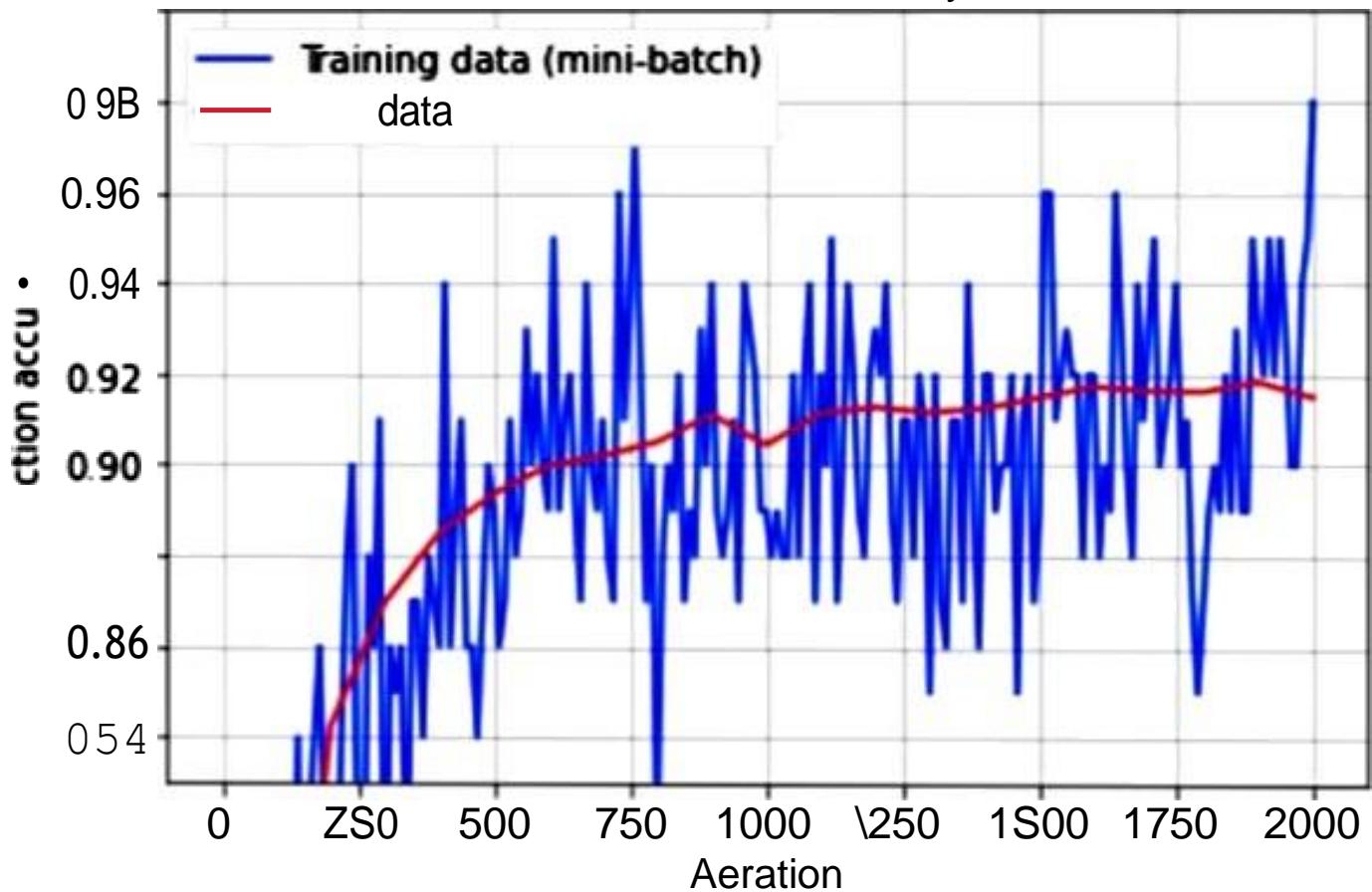
Plot

(v) Yes, using adjusted learning rate & with 200 & 100 layers. Accuracy is above 99%.

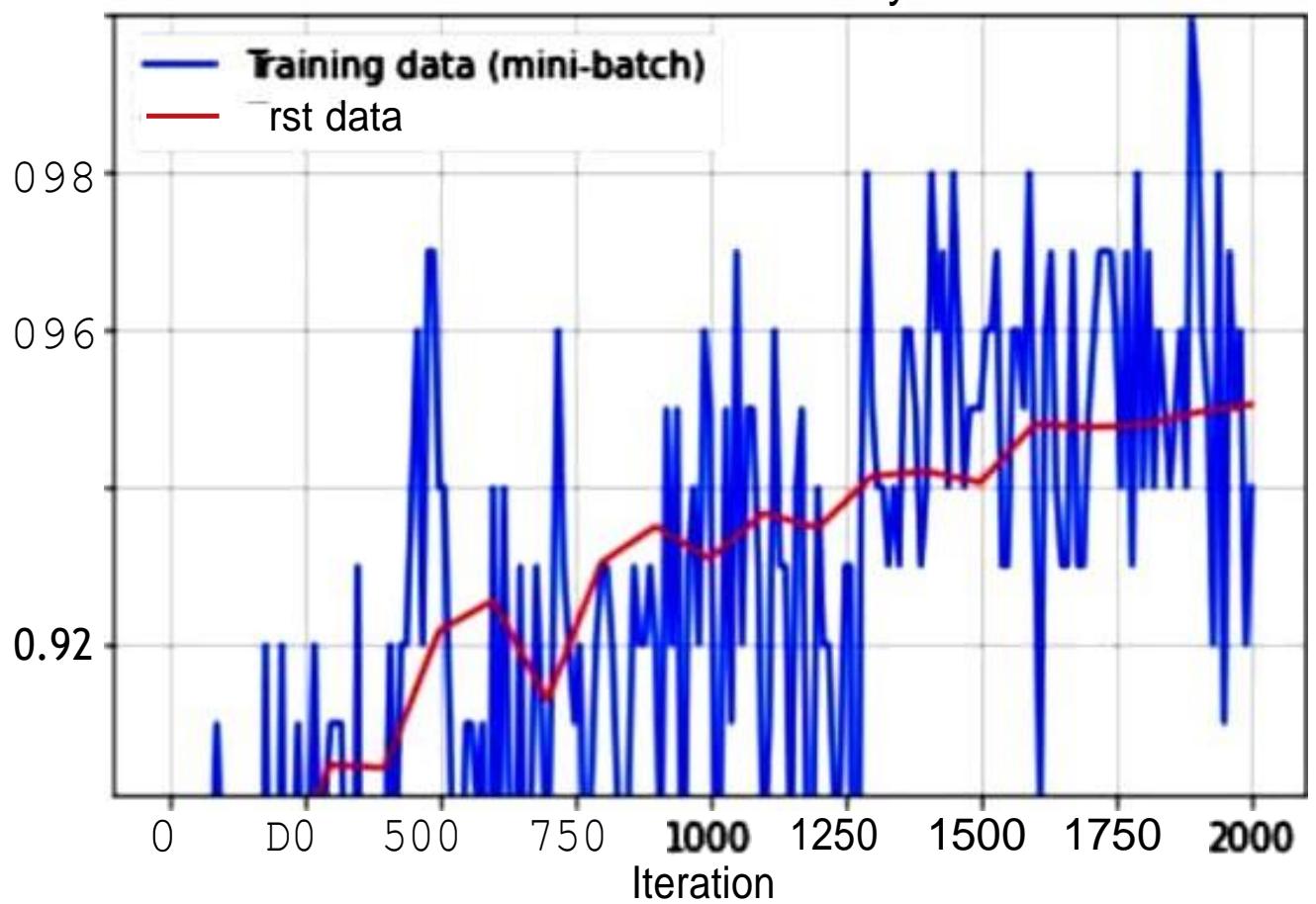
Prediction accuracy



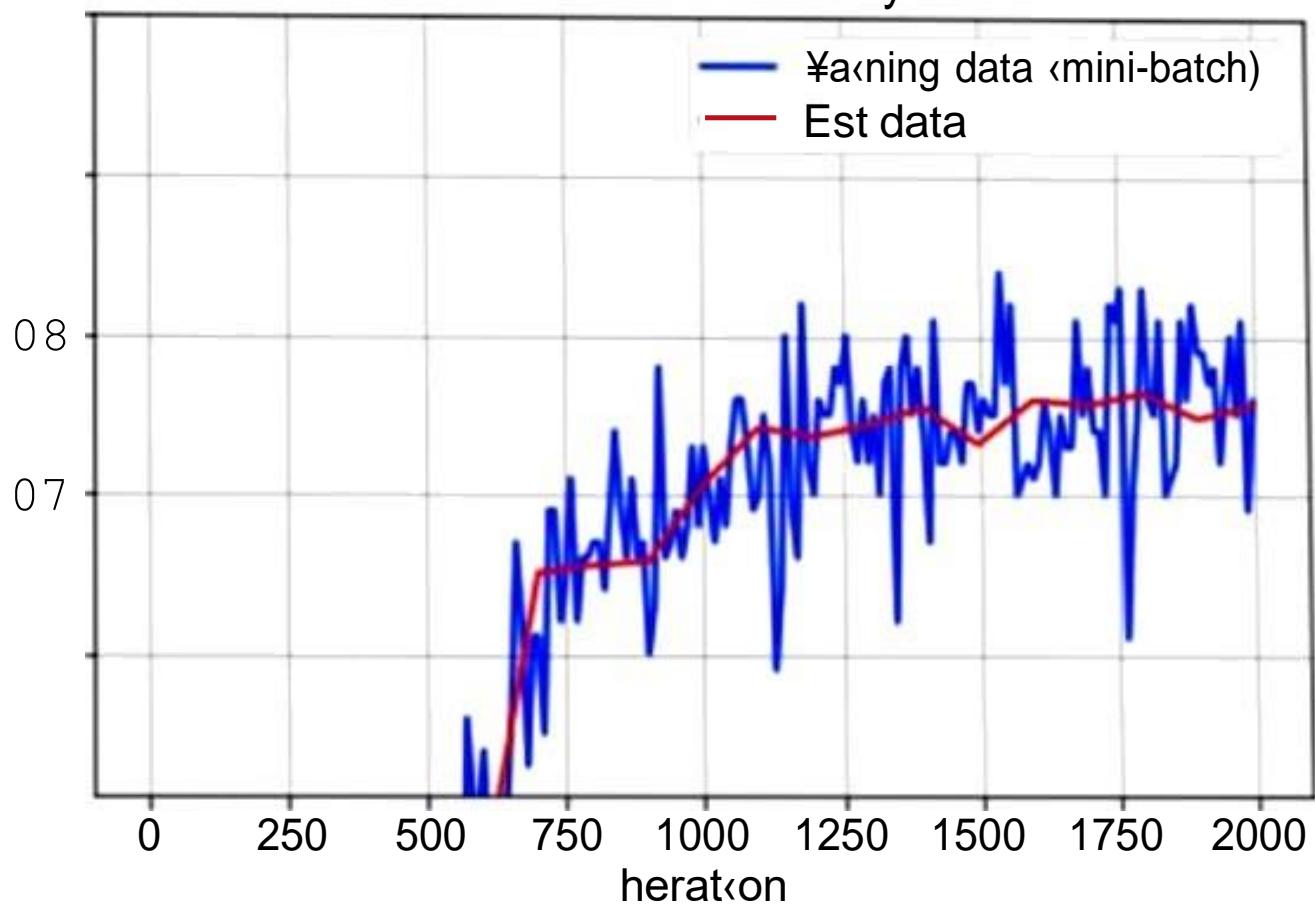
Prediction accuracy

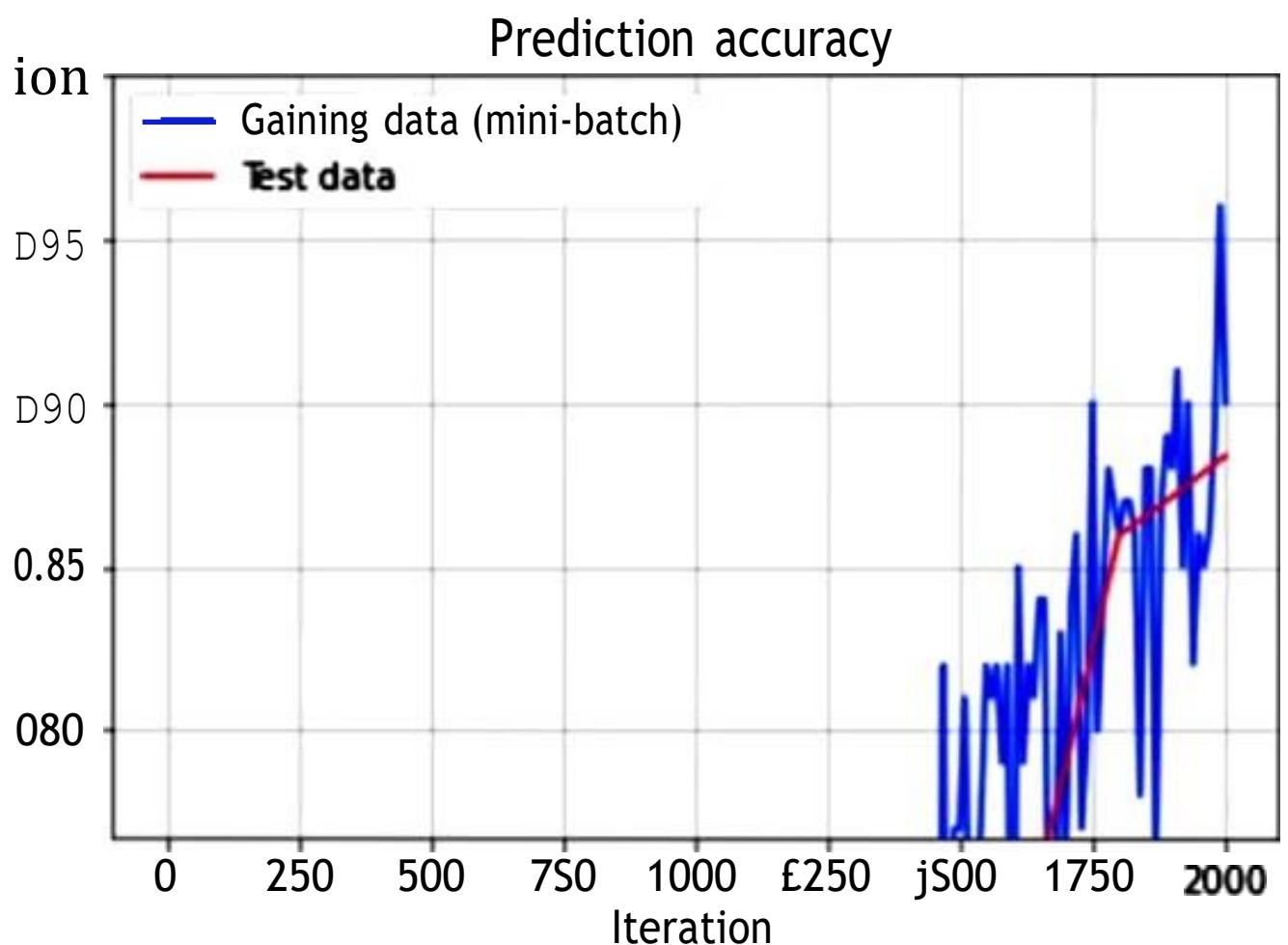


Prediction accuracy

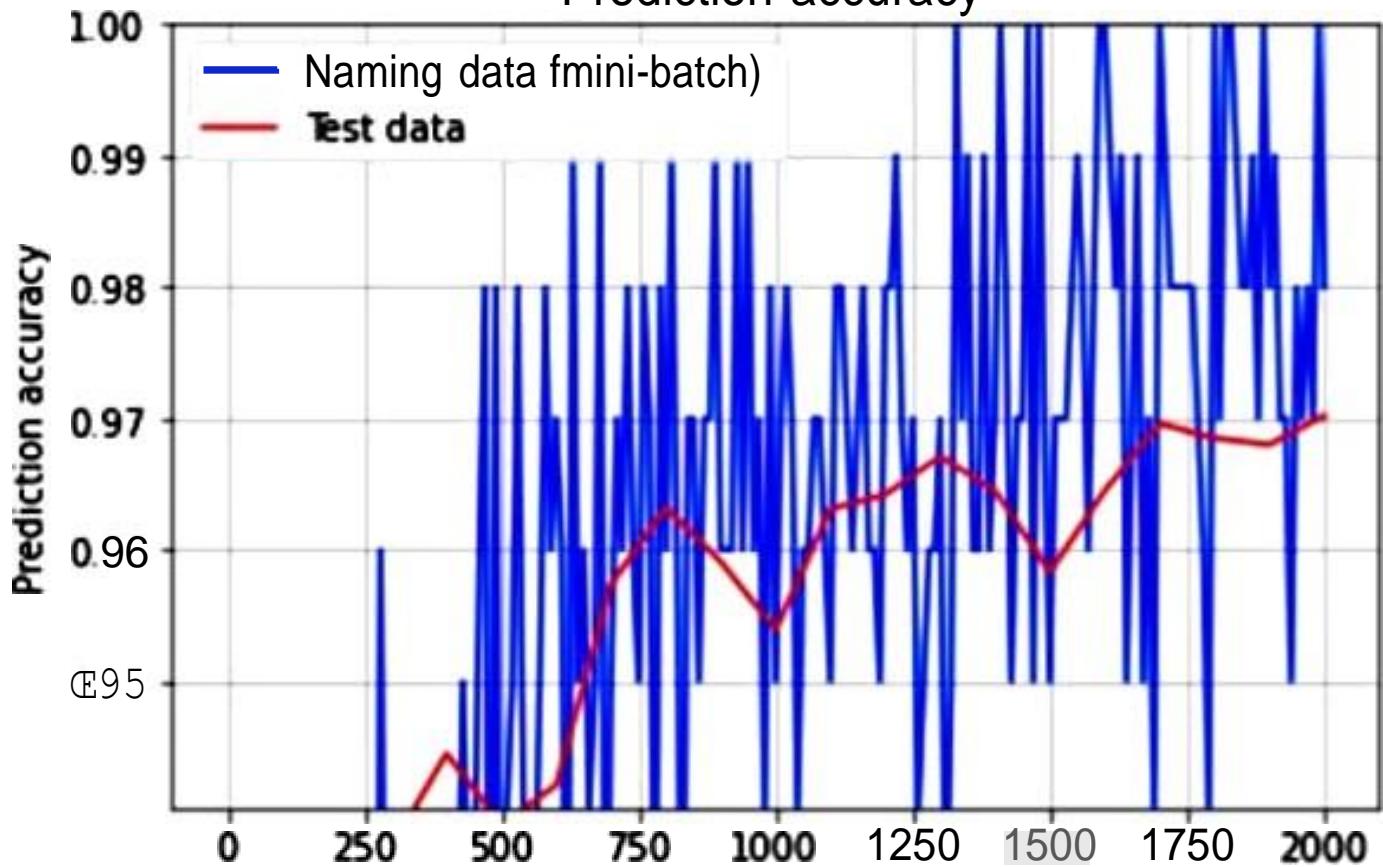


Prediction accuracy

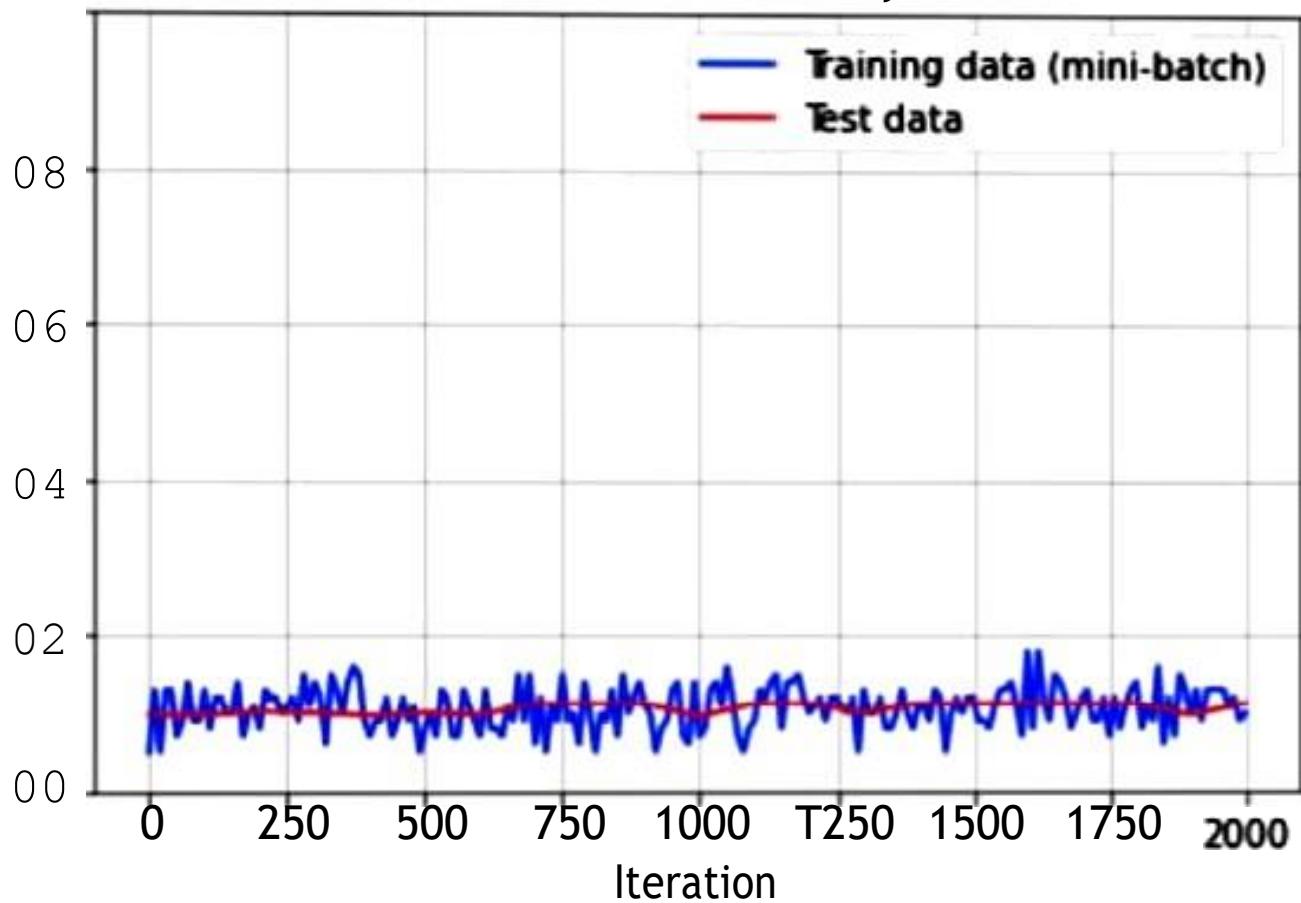




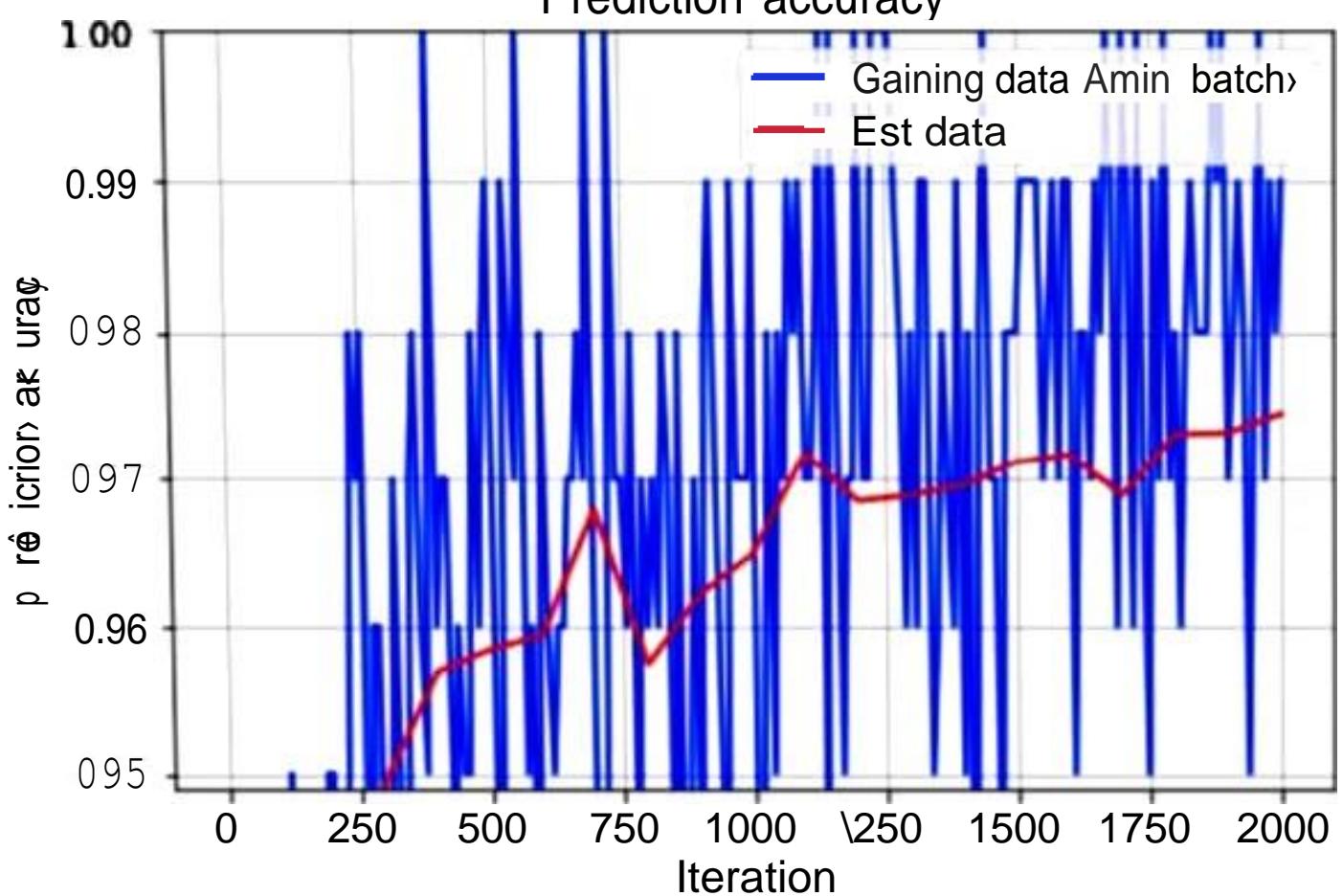
Prediction accuracy



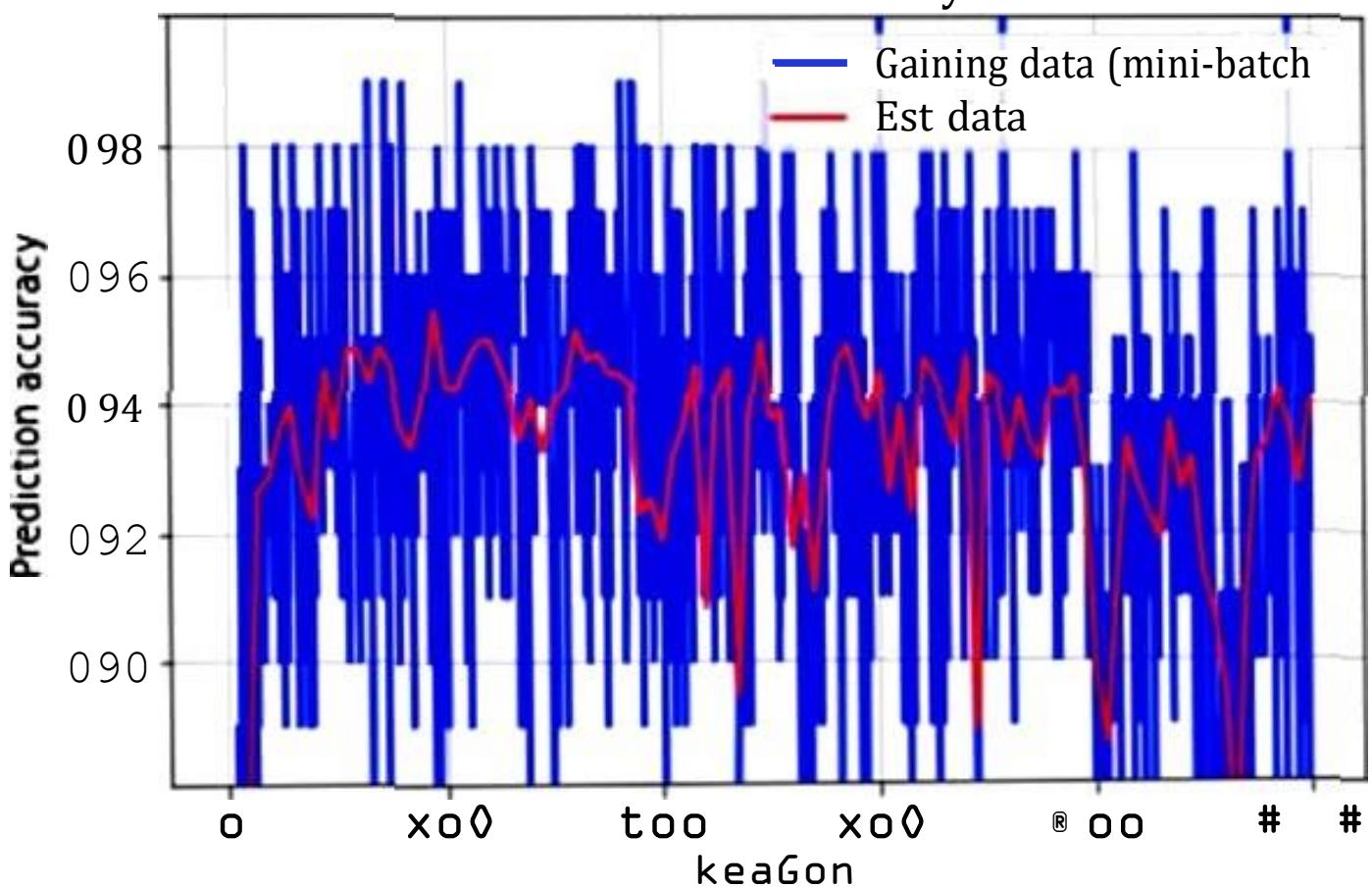
Prediction accuracy



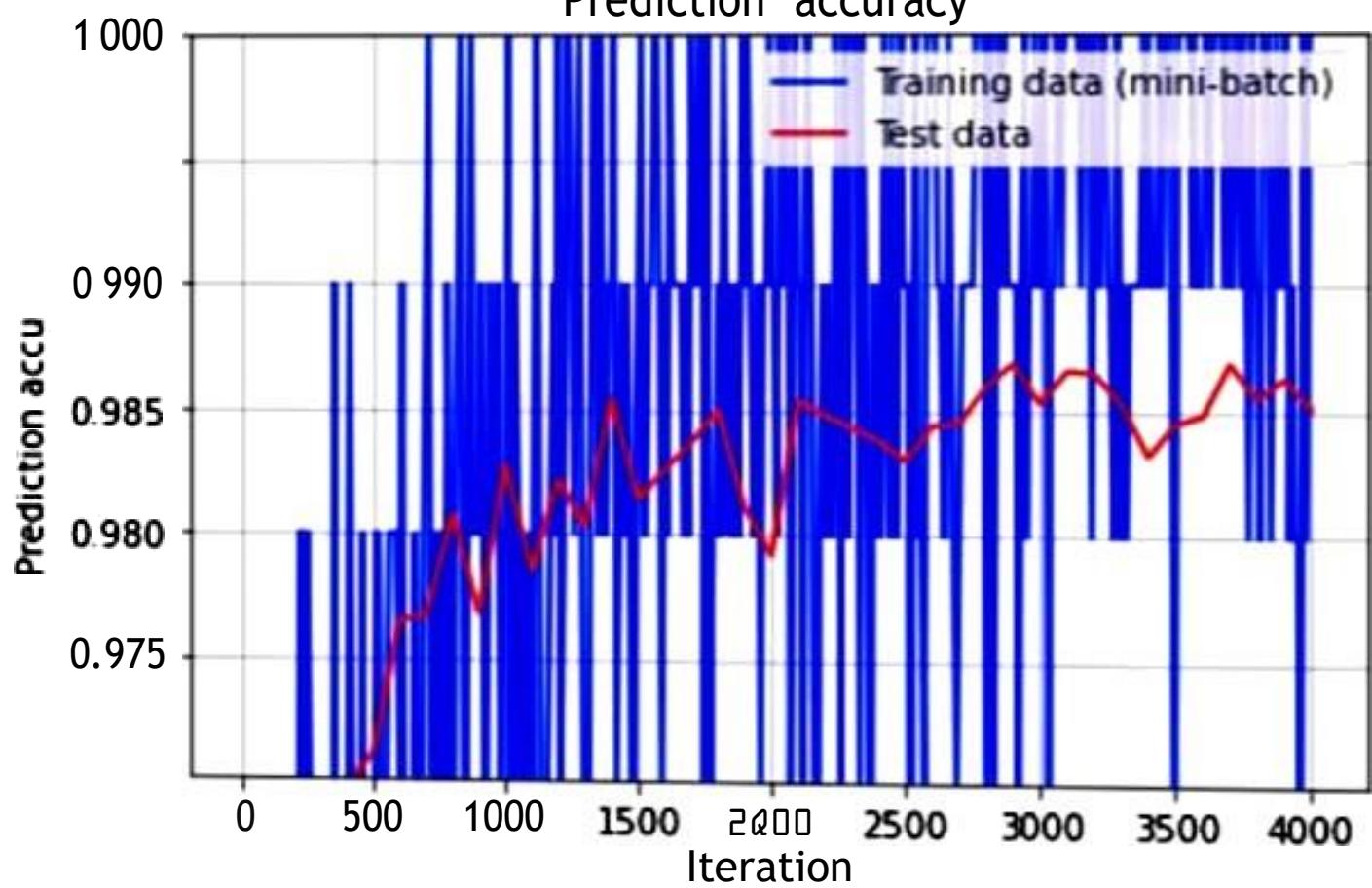
Prediction accuracy

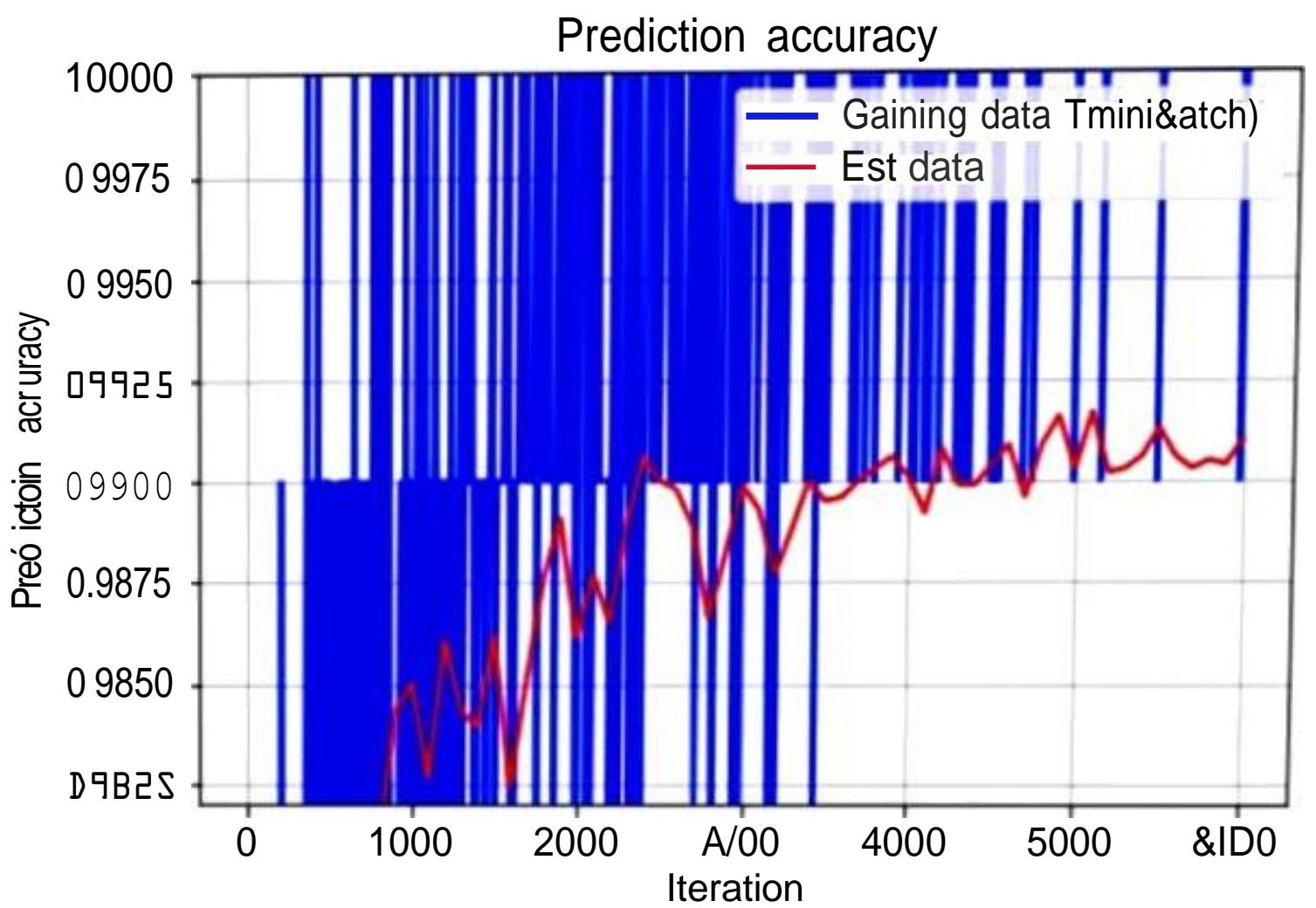


Prediction accuracy



Prediction accuracy





Prediction accuracy

