1.

2.



2.

PLA: $w_{t+1} = w_t + [y_t \neq \text{sign}(w_t^T x)] y x$

SGD: $w_{t+1} = w_t - \eta \nabla E_{in}(w_t)$.    $\eta$ : learning rate

$\begin{cases} y=+1,\ w_t^T x > 0 \begin{cases} y_t \neq \text{sign}(w_t^T x) \to \text{false} \Rightarrow w_{t+1} = w_t \\ \nabla E_{in} = 0 \Rightarrow w_{t+1} = w_t \end{cases} \\ y=-1,\ w_t^T x < 0 \begin{cases} y_t \neq \text{sign}(w_t^T x) \to \text{true} \Rightarrow w_{t+1} = w_t + x \\ \nabla E_{in} = \nabla err(w,x,y) = \nabla(-w^T x) = -x \Rightarrow w_{t+1} = w_t + x. \end{cases} \end{cases}$

$\begin{cases} y=-1,\ w_t^T x > 0 \begin{cases} y_t \neq \text{sign}(w_t^T x) \to \text{true} \Rightarrow w_{t+1} = w_t - x \\ \nabla E_{in} = \nabla err(w,x,y) = \nabla(w^T x) = x \Rightarrow w_{t+1} = w_t - x \end{cases} \\ y=-1,\ w_t^T x < 0 \begin{cases} y_t \neq \text{sign}(w_t^T x) \to \text{false} \Rightarrow w_{t+1} = w_t \\ \nabla E_{in} = 0 \Rightarrow w_{t+1} = w_t \end{cases} \end{cases}$

→ 由上得知, SGD using error function $\max(0, -y w^T x)$ results in PLA.

3. $\hat{E}_2(\Delta u, \Delta v) = b_{uu}(\Delta u)^2 + b_{vv}(\Delta v)^2 + b_{uv}(\Delta u)(\Delta v) + b_u \Delta u + b_v \Delta v + b$. $\approx E(u+\Delta u, v+\Delta v)$

target : $\nabla \hat{E}_2(\Delta u, \Delta v) = 0 \longleftrightarrow \nabla E(u+\Delta u, v+\Delta v) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$\nabla E(u+\Delta u, v+\Delta v) \approx \begin{bmatrix} b_{uu} \Delta u + b_{uv} \Delta v + b_u \\ b_{vv} \Delta v + b_{uv} \Delta u + b_v \end{bmatrix} = \begin{bmatrix} b_u \\ b_v \end{bmatrix} + \begin{bmatrix} b_{uu} & b_{uv} \\ b_{uv} & b_{vv} \end{bmatrix} \times \begin{bmatrix} \Delta u \\ \Delta v \end{bmatrix} = \nabla E(u,v) + \nabla^2 E(u,v) \begin{bmatrix} \Delta u \\ \Delta v \end{bmatrix}$

$\because \nabla^2 E(u,v)$ is positive define $\because \lambda > 0$, $\lambda \in \lambda(\nabla^2 E(u,v))$, $\because \lambda > 0, \lambda \in \lambda(\nabla^2 E(u,v))$) $\therefore \nabla^2 E(u,v)^{-1}$ exits

$\therefore \begin{bmatrix} \Delta u \\ \Delta v \end{bmatrix} = -(\nabla^2 E(u,v))^{-1} \nabla E(u,v)$

4. $y = \{1, 2, \dots, k\}$.   $h_y(x) = (\exp(w_y^T x)) / (\sum_{k=1}^{k} \exp(w_k^T x))$

→ maximum likelihood of $h_y(x)$

→ $\max_h \prod_{n=1}^{N} h_y(x_n) \to -\min \frac{1}{N} \ln \prod_{n=1}^{N} \exp(w_y^T x_n) - \ln \prod_{n=1}^{N} (\sum_{k=1}^{k} \exp(w_k^T x_n))$

→ $-\min \frac{1}{N} \sum_{n=1}^{N} w_y^T x_n - \sum_{n=1}^{N} \ln(\sum_{k=1}^{k} \exp(w_k^T x_n)) \Rightarrow \min \frac{1}{N} \sum_{n=1}^{N} (\ln(\sum_{k=1}^{k} \exp(w_k^T x_n)) - w_y^T x_n)$

$\therefore E_{in}(w_1, \dots, w_k) = \frac{1}{N} \sum_{n=1}^{N} \ln(\sum_{k=1}^{k} \exp(w_k^T x_n) - w_y^T x_n)$.

5. $x = \{x_1, \dots, x_n\}^T$. $y = \{y_1, \dots, y_n\}^T$. $\tilde{x} = \{\tilde{x}_1 \dots \tilde{x}_k\}^T$. $\tilde{y} = \{\tilde{y}_1 \dots \tilde{y}_k\}^T$.

$\min_w \frac{1}{N+k} (\sum_{n=1}^{N} (y_n - w^T x_n)^2 + \sum_{k=1}^{k} (\tilde{y}_k - w^T \tilde{x}_k)^2)$

$\Rightarrow \min_w \frac{1}{N+k} [(w^T x^T x w + 2 w^T x^T y + y^T y) + (w^T \tilde{x}^T \tilde{x} w + 2 w^T \tilde{x}^T \tilde{y} + \tilde{y}^T y)] = E_{in}(w)$

$\nabla E_{in}(w) = \frac{2}{N+k} (x^T x w - x^T y + \tilde{x}^T \tilde{x} w - \tilde{x}^T \tilde{y}) = 0 \Rightarrow (x^T x + \tilde{x}^T \tilde{x}) w = x^T y + \tilde{x}^T \tilde{y}$

$w = (x^T x + \tilde{x}^T \tilde{x})^{-1} (x^T y + \tilde{x}^T \tilde{y})$.

6.

$$6. \quad w_{reg} = \arg\min_w \frac{\lambda}{N} \|w\|^2 + \frac{1}{N} \|Xw - y\|^2$$

$$\longrightarrow \quad \frac{2\lambda}{N} w_{reg} + \frac{2}{N} X^T (Xw_{reg} - y) = 0$$
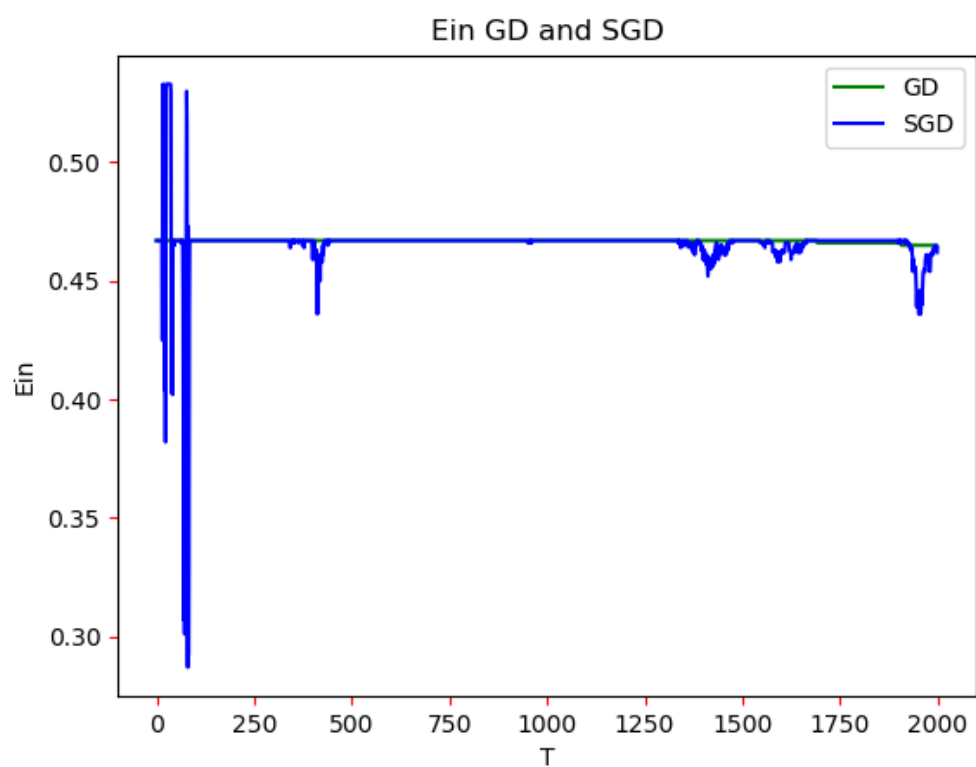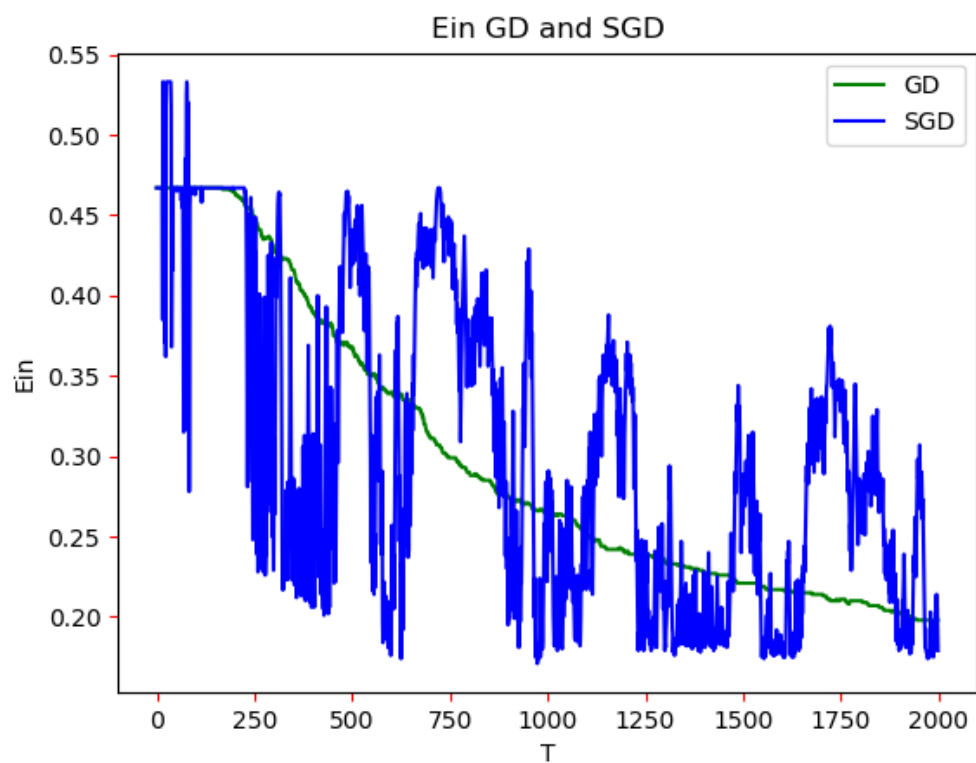
$$\longrightarrow \quad \frac{2\lambda}{N} w_{reg} + \frac{2}{N} (X^T X w_{reg} - X^T y) = 0$$

$$\longrightarrow \quad \lambda w_{reg} + X^T X w_{reg} - X^T y = 0$$

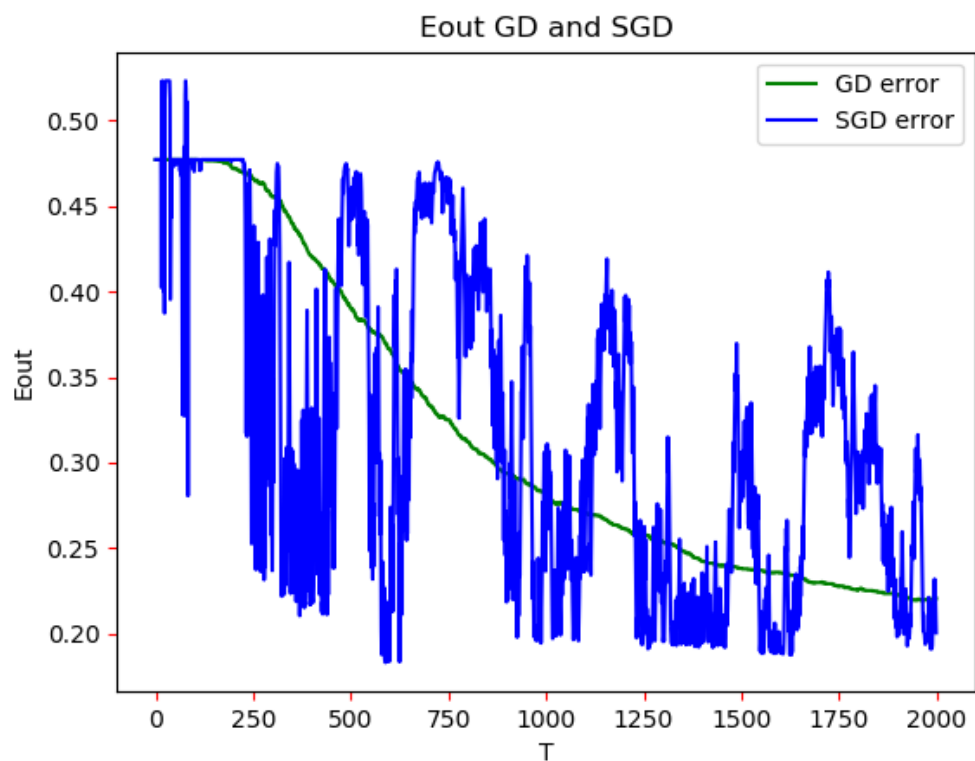$$\longrightarrow \quad (\lambda I + X^T X) w_{reg} - X^T y = 0$$

$$\longrightarrow \quad w_{reg} = (\lambda I + X^T X)^{-1} X^T y \quad \Rightarrow \quad w_{reg} = w \quad \Rightarrow \quad \hat{X} = \sqrt{\lambda} I, \quad \hat{y} = 0$$

7.

Learning rate 0.01 時，GD 與 SGD 的 error rate 有明顯的下降，但 Learning rate 0.001 時，error rate 變化不明顯，但 SGD 還是有機會達到 error 較低的時候。

8.

Learning rate 0.01 時，SGD 有較明顯的震盪，Learning rate 0.001 時，SGD 前幾次更新有較劇烈的震盪，但最後 error rate 趨近 GD。