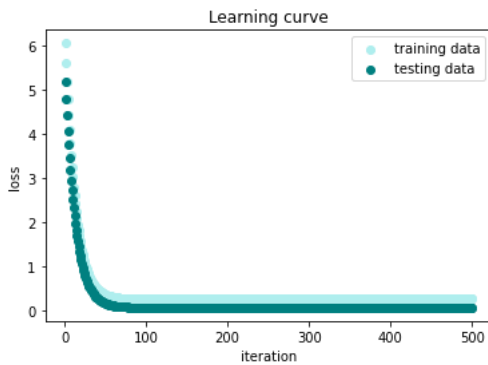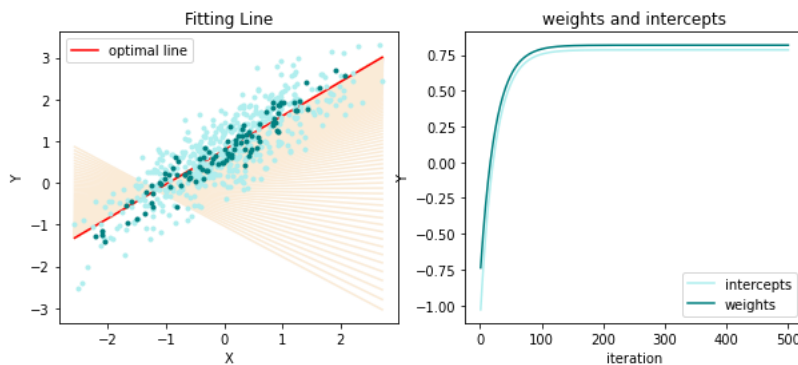1. Learning curve of the training :



2. MSE = 0.06870297339333518
3. Fitting Line : Y = 0.8179703746689468X + 0.7845650803186925 (weight ~= 0.818, intercept ~= 0.785)



4. What's the difference between Gradient Descent,(GD) Mini-Batch Gradient Descent(Mini-BGD), and Stochastic Gradient Descent(SGD)?

| | GD | Mini-BGD | SGD |
|---|---|---|---|
| Method | Parameters are updated after computing the gradient of error with respect to the entire dataset. | Parameters are updated after computing the gradient of error with respect to a subset of the dataset. | Parameters are updated after computing the gradient of error with respect to only a single data. |
| Noise | make smooth updates | Make less noisy updates than SGD (depend on the batch size) | make very noisy updates |
| Time | Take the most time | Balance between GD and SGD in terms of efficiency | Converges quickly for huge datasets. |
| convergence ratio | quick convergence ratio to a global minimum if the loss function is convex (and to local minimum one for non-convex functions) | Worse than GD | The worst (needs more iterations) |
| tolerance | | Can "bounce around" global optimum — may require the bigger tolerance. | Can "bounce around" global optimum — may require the biggest tolerance. |

1.

(1) $0.2\left(\frac{3}{10}\right) + 0.4\left(\frac{2}{4}\right) + 0.4\left(\frac{4}{20}\right) = 0.06 + 0.2 + 0.08$

$= 0.34$ ∎

(2) $P(🍎 \wedge B \mid 🍎)$

$= \dfrac{0.4(0.5)}{0.2(0.3) + 0.4(0.5) + 0.4(0.6)} = \dfrac{0.2}{0.5}$

$= \dfrac{2}{5}$ ∎

2. Let $R_1$ be the distribution area of class $C_1$, and $R_2$ be the area of class $C_2$, we know

① $P(\text{mistake}) = \displaystyle\int_{R_1} P(x, c_2)\,dx + \int_{R_2} P(x, c_1)\,dx$

② In the error made in $R_1$:

we always have $P(c_1|x) \geq P(c_2|x)$,

which implies $P(c_2|x) \leq \{P(c_1|x)\,P(c_2|x)\}^{\frac{1}{2}}$

$\displaystyle\int_{R_1} P(x, c_2)\,dx = \int_{R_1} P(c_2|x)\,P(x)$

$\leq \displaystyle\int_{R_1} \{P(c_1|x)\,P(c_2|x)\}^{\frac{1}{2}} P(x)\,dx$

$= \displaystyle\int_{R_1} \{P(x,c_1)\,P(x,c_2)\}^{\frac{1}{2}}\,dx$

③ In the error made in $R_2$:

we always have $P(c_2|x) \geq P(c_1|x)$,

which implies $P(c_1|x) \leq \{P(c_1|x)\,P(c_2|x)\}^{\frac{1}{2}}$

$\displaystyle\int_{R_2} P(x, c_1)\,dx = \int_{R_2} P(c_1|x)\,P(x)\,dx$

$\leq \displaystyle\int_{R_2} \{P(c_1|x)\,P(c_2|x)\}^{\frac{1}{2}} P(x)\,dx$

$= \displaystyle\int_{R_2} \{P(x,c_1)\,P(x,c_2)\}^{\frac{1}{2}}\,dx$

Substitute ②③ back to ①, we have

$P(\text{mistake}) = \displaystyle\int_{R_1} P(x, c_2)\,dx + \int_{R_2} P(x, c_1)\,dx$

$\leq \displaystyle\int_{R_1} \{P(x,c_1)\,P(x,c_2)\}^{\frac{1}{2}}\,dx + \int_{R_2} \{P(x,c_1)\,P(x,c_2)\}^{\frac{1}{2}}\,dx$

$= \displaystyle\int \{P(x,c_1)\,P(x,c_2)\}^{\frac{1}{2}}\,dx$ ∎

3.

(1) ① Assume $g(Y) = E[X|Y]$

$\Rightarrow g(Y) = \displaystyle\sum_{x \in X} E[X = x | Y]\,P(x)$

$\Rightarrow g(Y) = E_x[x|Y]$

② $E[X] = \displaystyle\sum_{y \in Y} E[X|Y = y]\,P(y)$

$= \displaystyle\sum_{y \in Y} g(y)\,P(y)$

$= E_y[g(y)]$

$= E_y[E_x[x|Y]]$ ∎

(2) ① $\mathrm{Var}_x(x|Y) = E_x[x^2|Y] - (E_x[x|Y])^2$

$\Rightarrow E_y[\mathrm{Var}_x(x|y)] = E_y[E_x[x^2|Y]] - E_y[(E_x[x|Y])^2]$

$= E[x^2] - E_y[(E_x[x|Y])^2]$

② $\mathrm{Var}_y(E_x[x|Y]) = E_y[(E_x[x|Y])^2] - (E_y[E_x[x|Y]])^2$

①② $\Rightarrow E_y[\mathrm{Var}_x(x|y)] + \mathrm{Var}_y(E_x[x|Y])$

$= E[x^2] - \cancel{E_y[(E_x[x|Y])^2]} + \cancel{E_y[(E_x[x|Y])^2]} - \underline{(E_y[E_x[x|Y]])^2}_{(E[x])^2}$

$= E[x^2] - (E[x])^2 = \mathrm{Var}(x)$ ∎