1.
$$f_{0/1}(x) = \arg\max_{y \in \{-1,1\}} P(y|x) = \text{sign}\left(P(1|x) - \frac{1}{2}\right)$$

| $f_{0/1}$ | $g$ +1 | $g$ -1 | |
|---|---|---|---|
| +1 | no err | 1 | |
| -1 | 1 | no err | |

For CIA case,

| $f_{CIA}$ | $g$ +1 | $g$ -1 | |
|---|---|---|---|
| +1 | no err | 1 | |
| -1 | 1000 | no err | |

$$\begin{cases} 1 & , y_n \neq g(x_n) \wedge y_n = +1 \\ 1000 & , y_n \neq g(x_n) \wedge y_n = -1 \end{cases}$$

Let $f_{CIA} = \text{sign}(P(1|x) - \alpha)$

Expected cost classifying $x$ as positive : $1000\, P(-1|x)$

Expected cost classifying $x$ as negative : $1 \cdot P(1|x)$

We only classify $x$ as positive when the cost that we classify $x$ as negative is less than classifying as positive.

$$\therefore\ 1000\, P(-1|x) < 1 \cdot P(1|x)$$

By $1 - P(1|x) = P(-1|x)$

$$1000\,(1 - P(1|x)) < P(1|x)$$

$$P(1|x) > \frac{1000}{1001} \Rightarrow \alpha = \frac{1000}{1001}$$

$$\therefore f_{CIA} = \text{sign}\left(P(1|x) - \frac{1000}{1001}\right)$$

#

#

2. $P(y = +f(x) | x) = 1 - \varepsilon$

$P(y = -f(x) | x) = \varepsilon$

Given $E_{out}(g) = \mathbb{E}_{x \sim P(x)} [[ g(x) \neq f(x) ]]$.

Find $\mathbb{E}_{(x,y) \sim P(x,y)} [[ g(x) \neq y ]]$

<u>case 1.</u> when $g(x) = f(x)$ but $y \neq f(x)$ (due to noise)

$\quad E_1 = (1 - E_{out}(g)) \cdot (\varepsilon) \longrightarrow ①$

<u>case 2.</u> when $g(x) \neq f(x)$, but

$\quad\quad\quad\quad y(x) = f(x)$

$\quad E_2 = E_{out}(g) \cdot (1 - \varepsilon)$

|  | $[[ y(x) = f(x) ]]$ | |
|---|---|---|
|  | $+1$ | $-1$ |
| $[[ g(x) = f(x) ]] \rightarrow +1$ | $\times$ | err |
| $-1$ | err | $\times$ |

$\therefore E_{(x,y) \sim P(x,y)} [[ g(x) \neq y ]] = (1 - E_{out}(g)) \varepsilon + E_{out}(g)(1-\varepsilon)$

$$= \varepsilon - 2 E_{out}(g) + E_{out}(g)$$

3. $h(x) = wx, \quad E_{in}(w) = \frac{1}{N} \sum_{n=1}^{N} (h(x_n) - y_n)^2 = \frac{1}{N} \sum_{n=1}^{N} (wx_n - y_n)^2$ #

$$= \frac{1}{N} \sum_{n=1}^{N} (w^2 x_n^2 - 2wx_n y_n + y_n^2)$$

To find $\min_{w} E_{in}(w) \Rightarrow$ We take $\nabla E_{in}(w) = 0$

$\therefore \nabla E_{in}(w) = \frac{\partial E}{\partial w} = \frac{1}{N} \sum_{n=1}^{N} (2x_n^2 w - 2x_n y_n) = 0$

$\Rightarrow \sum_{n=1}^{N} 2x_n^2 w - \sum_{n=1}^{N} 2x_n y_n = 0 \Rightarrow W_{lin} = \frac{\sum_{n=1}^{N} x_n y_n}{\sum_{n=1}^{N} x_n^2}$ #

4.

$f(x) = ax^2 + b$ , $h(x) = w_0 + w_1 x$ , $x$ sampled from $[0, 1]$

$$\therefore E_{sqr}(w_0, w_1) = \int_0^1 (h(x) - f(x))^2 dx$$

$$= \int_0^1 \left((w_0 + w_1 x) - (ax^2 + b)\right)^2 dx = \int_0^1 \left(-ax^2 + w_1 x + w_0 - b\right)^2 dx$$

To Find $\min\limits_{w} E_{sqr}(w_0, w_1) \Rightarrow \nabla E(w_0, w_1) = 0$

$$\frac{\partial E}{\partial w_0} = 0 \Rightarrow 2\int_0^1 (-ax^2 + w_1 x + w_0 - b) \cdot 1\, dx = 0$$

$$\Rightarrow 2\left(\frac{-a}{3}x^3 + \frac{w_1}{2}x^2 + (w_0 - b)x\right)\Big|_0^1 = 0$$

$$\Rightarrow \frac{-a}{3} + \frac{w_1}{2} + (w_0 - b) = 0 \Rightarrow w_1 + 2w_0 = \frac{2}{3}a + 2b \qquad ——①$$

$$\frac{\partial E}{\partial w_1} = 0 \Rightarrow 2\int_0^1 (-ax^2 + w_1 x + w_0 - b) \cdot (x)\, dx = 0$$

$$\Rightarrow 2\int_0^1 (-ax^3 + w_1 x^2 + (w_0 - b)x)\, dx = 0$$

$$\Rightarrow 2\left(\frac{-a}{4}x^4 + \frac{w_1}{3}x^3 + \frac{w_0 - b}{2}x^2\right)\Big|_0^1 = 0$$

$$\Rightarrow \frac{-a}{4} + \frac{w_1}{3} + \frac{w_0 - b}{2} = 0 \Rightarrow 4w_1 + 6(w_0 - b) - 3a = 0$$

$$\Rightarrow 4w_1 + 6w_0 = 3a + 6b \qquad ——②$$

By ①②,

$$\therefore w_1^* = a, \quad w_0^* = \frac{-a}{6} + b$$

$$\therefore (w_0^*, w_1^*) = \left(\frac{-a}{6} + b, a\right) \quad \text{We have } \min_{w} E_{in}(w_0, w_1) \#$$

5.

$$W_{Lin} = (X^T X)^{-1} X^T y$$

$$E(W'_{Lin}) = \frac{1}{N} \left\| X W'_{Lin} - \left( ay + \begin{bmatrix} b \\ b \\ \vdots \\ b \end{bmatrix} \right) \right\|^2$$

$$= \frac{1}{N} \left( W'^T_{Lin} X^T X W'_{Lin} - 2 \left( X W'_{Lin} \left( ay + \begin{bmatrix} b \\ b \\ \vdots \\ b \end{bmatrix} \right) \right) + \left( ay + \begin{bmatrix} b \\ b \\ \vdots \\ b \end{bmatrix} \right)^2 \right)$$

To find $\min_W E_{in}(W'_{Lin})$

$$\nabla E(W'_{Lin}) = 0$$

$$\Rightarrow 2 X^T X W'_{Lin} - 2 X^T \left( ay + \begin{bmatrix} b \\ b \\ \vdots \\ b \end{bmatrix} \right) = 0$$

Let $W'_{Lin} = W_1 + W_2$ (By linear combination)

For $W_1$, $X^T X W_1 - X^T ay = 0 \Rightarrow W_1 = a (X^T X)^{-1} X^T y$

$$= a W_{Lin}$$

For $W_2$, $X^T X W_2 - X^T \begin{bmatrix} b \\ b \\ \vdots \\ b \end{bmatrix} = X^T \left( X W_2 - \begin{bmatrix} b \\ b \\ \vdots \\ b \end{bmatrix} \right) = 0$

By $X = \begin{bmatrix} | & - x_1^T - \\ | & - x_2^T - \\ | & \vdots \\ | & - x_N^T - \end{bmatrix}$ $\therefore W_2 = \begin{bmatrix} b \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

$$\therefore W_{Lin} = a W_{Lin} + \begin{bmatrix} b \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \left( \begin{bmatrix} b \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ is } (d+1) \times 1 \text{ column vector} \right)$$

\#

6.

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^{N} ln\left(1 + exp(-y_n w^T x_n)\right) , \quad h_t = \frac{1}{1 + exp(w_t^T x)} \quad \text{Given}$$

$$\nabla E_{in}(w) = \frac{\partial E}{\partial w} = \frac{1}{N} \sum_{n=1}^{N} (-y_n^t x_n)\left(\frac{exp(-y_n w^T x_n)}{1 + exp(-y_n w^T x_n)}\right) \quad \begin{pmatrix} \because h(s) \\ = 1 - h(s) \end{pmatrix}$$

$$= \frac{1}{N} \sum_{n=1}^{N} (-y_n x_n) \, h_t(y_n w^T x_n) \quad \begin{array}{l} \text{by definition of} \\ \text{logistic function} \end{array}$$

For $\nabla^2 E_{in}(w) = \frac{\partial}{\partial w_j} \frac{\partial}{\partial w_i}(E_{in}(w))$ , We substitute $k = exp(-y_n w^T x_n)$

$$= \left(\frac{1}{N} \sum_{n=1}^{N} (-y_{ni} x_{ni}) \frac{(-y_{nj} x_{nj}) k \cdot (1+k) - k(-y_{nj}x_{nj})k}{(1+k)^2}\right)$$

$$\underbrace{=}_{(i,j)} \frac{1}{N} \sum_{n=1}^{N} \left(x_{ni} \, x_{nj} \, y_{ni} \, y_{nj} \, h_t(y_n w_t^T x_n)(1 - h_t(y_n w_t^T x_n))\right)$$

The single term denote Hessian Matrix

⇒ Express the sum in matrix form :

Let $X$ as matrix is a data point $x_n$, $D_{nn}$ is a diagonal matrix
where $n^{th}$ diagonal entry :

$$\therefore A_E(w_t)\bigg|_{E = E_{in}} = X^T D X$$

$$\therefore D_{nn} = y_n^2 \, h_t(y_n w_t^T x_n)\left(1 - h_t(y_n w_t^T x_n)\right)$$

$\therefore$ For diagonal Matrix $D$ ,

$$D = diag\left[ y_1^2 h_t(y_1 w_t^T x_1)(1 - h_t(y_1 w_t^T x_1)), \; y_2^2 h_t(y_2 w_t^T x_2)(1 - h_t(y_2 w_t^T x_2)) \right.$$
$$\left. \cdots , \; y_n^2 h_t(y_n w_t^T x_n)(1 - h_t(y_n w_t^T x_n)) \right]$$

7. Given $err(s,y) = (max(0, 1-ys))^2$, $s = w^Tx$

By SGD: $w_{t+1} \leftarrow w_t + \eta(-\nabla err(w, x_n, y_n))$ ($\eta$ is fixed learning rate)

<u>case1.</u> if $1+ys \geq 0 \Rightarrow err(s,y) = (1-ys)^2$

$\therefore \dfrac{\partial err(s,y)}{\partial w} = \dfrac{\partial err(s,y)}{\partial s} \cdot \dfrac{\partial s}{\partial w}$

$= -2(1-ys)y \cdot x = -2yx(1-ys)$

<u>case2.</u> if $1-ys \leq 0 \Rightarrow err(s,y) = 0 \Rightarrow$ The $\nabla err(s,y) = 0$

$\therefore$ SGD update: $w_{t+1} \leftarrow w_t + \eta(2yx(1-ys))$

(only when $1-ys \geq 0$) #

- Comparing to original PLA, Explanations:

① Loss Function: original PLA is 0/1 loss, update only when misclassification, while the new approach uses the truncated squared loss which smoothen the loss surface.

② Update rule: In original PLA $w_t \leftarrow w_t + \eta yx$ for misclassified points. The truncated squared loss SGD updates with weighted factor scaled by $2(1-ys)$

(Furthermore, only updates when $1-ys \geq 0$)

③ Convergence: The original PLA might not converge for non-linearly separable data. The truncated squared SGD might converge to a solution even if the data is not linearly-separable due to nature of truncated squared loss #

8.

$$h_y(x) = \frac{\exp(W_y^T x)}{\sum_{i=1}^{k} \exp(W_i^T x)}$$

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^{N} err(W, x_n, y_n)$$

$$= \frac{1}{N} \sum_{n=1}^{N} \sum_{k=1}^{K} [[y=k]] (-\ln h_k(x))$$

For single point $(x,y)$ (i.e. $y=k$)

$$err(W, x, y) = -\ln \frac{\exp(W_y^T x)}{\sum_{i=1}^{k} \exp(W_i^T x)} = -W_y^T x + \ln\left(\sum_{i=1}^{k} \exp(W_i^T x)\right)$$

① For $y=k$,

$$\nabla_{y=k}(err(w, x, y)) = \frac{\partial(-W_y^T x)}{\partial W_k} + \frac{\partial\left(\ln \sum_{i=1}^{k} \exp(W_i^T x)\right)}{\partial W_k}$$

$$= -x + \left. \frac{\exp(W_k^T x) \cdot x}{\sum_{i=1}^{k} \exp(W_i^T x)} \right|_{y=k}$$

$$= -x + h_k(x) \cdot x$$

② For $y \neq k$

$$\nabla_{y \neq k}(err(w, x, y)) = \frac{\partial\left(\ln \sum_{i=1}^{k} \exp(W_i^T x)\right)}{\partial W_k}$$

$$= \frac{\exp(W_k^T x) \cdot x}{\sum_{i=1}^{k} \exp(W_i^T x)} = h_k(x) \cdot x$$

$$\therefore \nabla err(w, x, y) = \begin{cases} -x + h_k(x) \cdot x & , \text{if } y=k \\ h_k(x) \cdot x & , \text{if } y \neq k \end{cases}$$

Double A

$$\therefore \nabla E_{in} = \frac{1}{N} \sum_{n=1}^{N} \nabla err(w, x_n, y_n)$$

$$\therefore \nabla E_{in} = \begin{cases} \frac{1}{N} \sum_{n=1}^{N} (-x + h_k(x) \cdot x) & , \text{ if } y = k \\ \frac{1}{N} \sum_{n=1}^{N} (h_k(x) \cdot x) & , \text{ if } y \neq k. \end{cases}$$

(Note each $\nabla err(w, x, y)$ is a matrix of size $(d+1) \times K$ #,

each column is the gradient of corresponding $W_k$ as derived

above.)

13. $$E_{in}(w) = \frac{1}{N} \sum_{n=1}^{N} \ln(1 + \exp(-y_n w_t^T x_n))$$

$$v = -(X^T D X)^{-1} \nabla E_{in}(w_t)$$

$$D_{nn} = h(x_n)(1 - h(x_n))$$

The analogy between logistic regression and linear regression

is that $$W_{Lin} = \underbrace{(X^T X)^{-1} X^T y}_{} = \underbrace{(\tilde{x}^T \tilde{x})^{-1}}_{\hookrightarrow (X^T D X)^{-1}} \underbrace{(\tilde{x}^T \tilde{y})}_{\hookrightarrow (\nabla E_{in}(w_t))}$$

To find analogy,

$$-\nabla E_{in}(w) = \frac{1}{N} \sum_{n=1}^{N} h_t(y_n W^T x_n)(+y_n x_n) = \frac{1}{N} \sum_{n=1}^{N} \frac{y_n x_n}{1 + \exp(y_n w^T x_n)}$$

(rewrite the sum)

$$\stackrel{\downarrow}{=} \frac{1}{N} \tilde{x}^T \frac{y_n}{1 + e^{y_n w^T x_n}}$$

Double A

$$\therefore \text{We define } \tilde{y} = \frac{1}{N\sqrt{D}} \left( \frac{y_n}{1 + \exp(y_n w^T x_n)} \right)$$

$$\tilde{x} = x\sqrt{D}$$

#

Verify: $\tilde{x}^T \tilde{x} = x^T \sqrt{D} \cdot \sqrt{D} x = x^T D x$ (correct)

$$\tilde{x}^T \tilde{y} = x^T \sqrt{D} \cdot \frac{1}{N\sqrt{D}} \left( \frac{y_n}{1 + \exp(y_n w^T x_n)} \right)$$

$$= \frac{1}{N} \left( x^T \cdot \left( \frac{y_n}{1 + \exp(y_n w^T x_n)} \right) \right)$$

$$= \frac{1}{N} \sum_{n=1}^{N} \left( \frac{1}{1 + \exp(y_n w^T x_n)} \right) (y_n x_n)$$

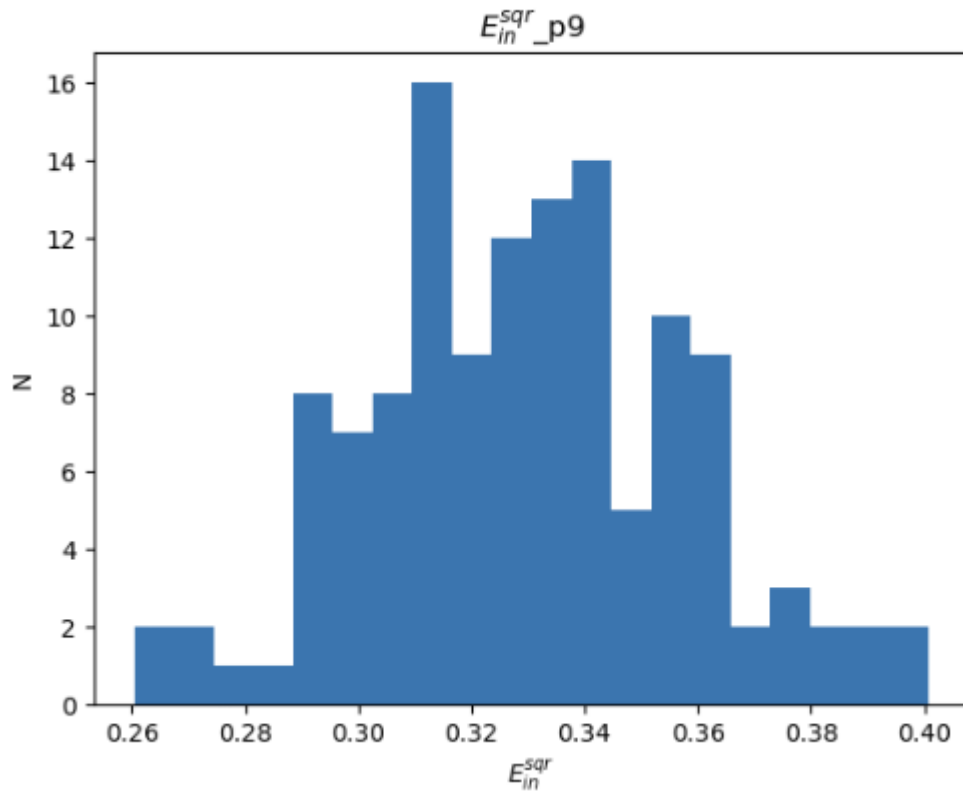$$= -\nabla E_{in}(w) \quad (\text{correct})$$
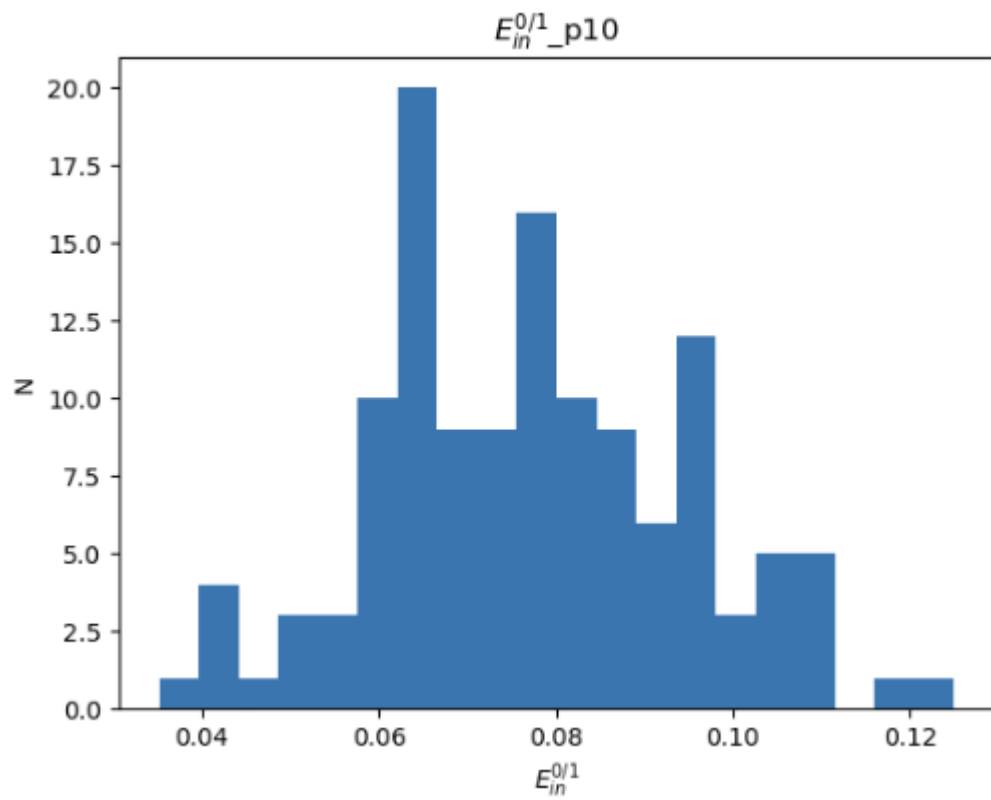
#

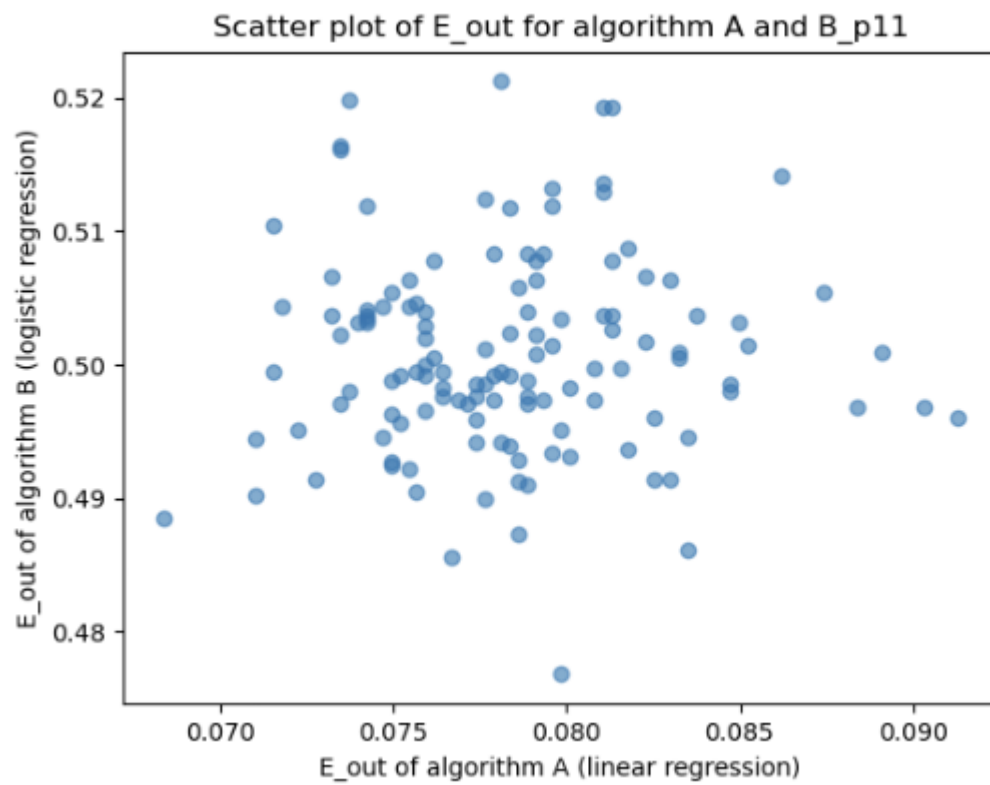# HTML hw3 solution

姓名: 謝銘倫, 系級:電機三, 學號:B10502166
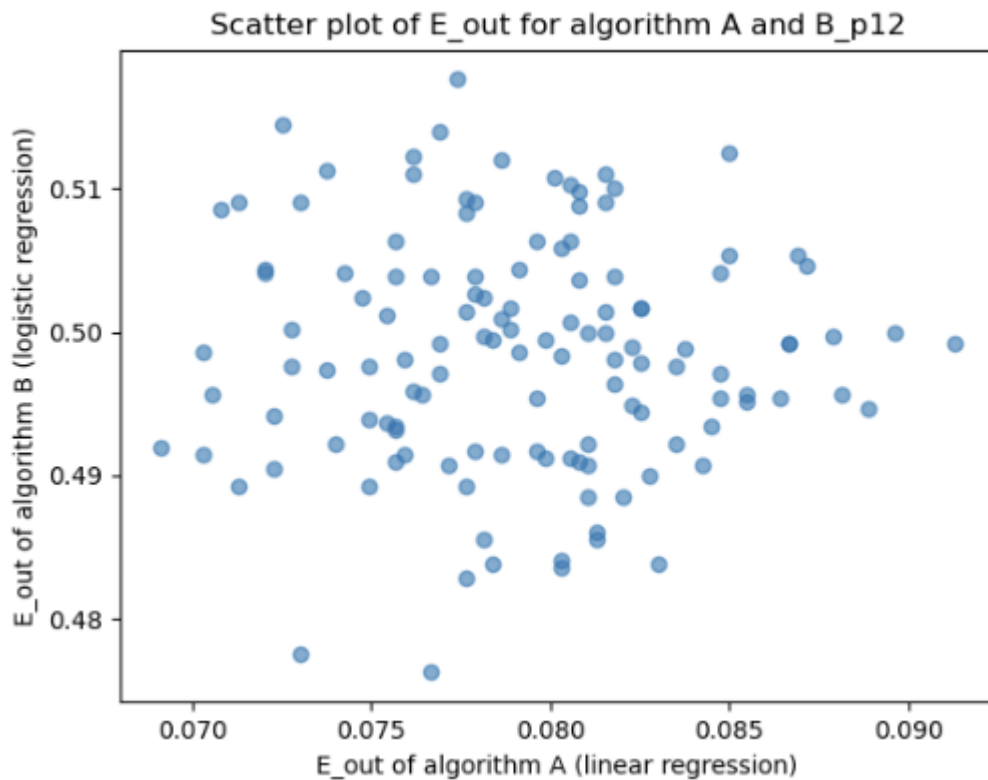
9.



median of E_in: 0.329

10.



median of E_in: 0.078

11.



Scatter plot of E_out for algorithm A and B_p11

median of $E_{out}(A(D))$ = 0.078

median of $E_{out}(B(D))$ = 0.499

12.



Scatter plot of E_out for algorithm A and B_p12

median of $E_{out}(A(D'))$ = 0.079

median of $E_{out}(B(D'))$ = 0.498

**Findings:**
**1.** 第11題與第12題的圖結果相似 median of E_out也相似

**2.** 因為logistic regression是用sigmoid function 這個函數是個monotone function並不容易受到outlier影響 因此結果相似

**3.** outlier的數量並不多 僅16筆 相對於原本training data的256筆 僅佔少數 因此對整體圖形的影響並不大