

1. (a)

Max Iterations	$E_{in}$	Binary $E_{in}$	Binary $E_{out}$	Time (s)
$10^4$	0.586	0.309	0.317	0.24
$10^5$	0.469	0.224	0.207	2.18
$10^6$	0.436	0.151	0.131	21.48

The logistic regression model generalizes well. The  $E_{in}$  and  $E_{out}$  are consistently very similar indicating that there is enough data for learning to occur, at least with this simple model. Increasing the maximum number of iterations, at least to up to  $10^6$ , continues to improve the model's performance.

(b)

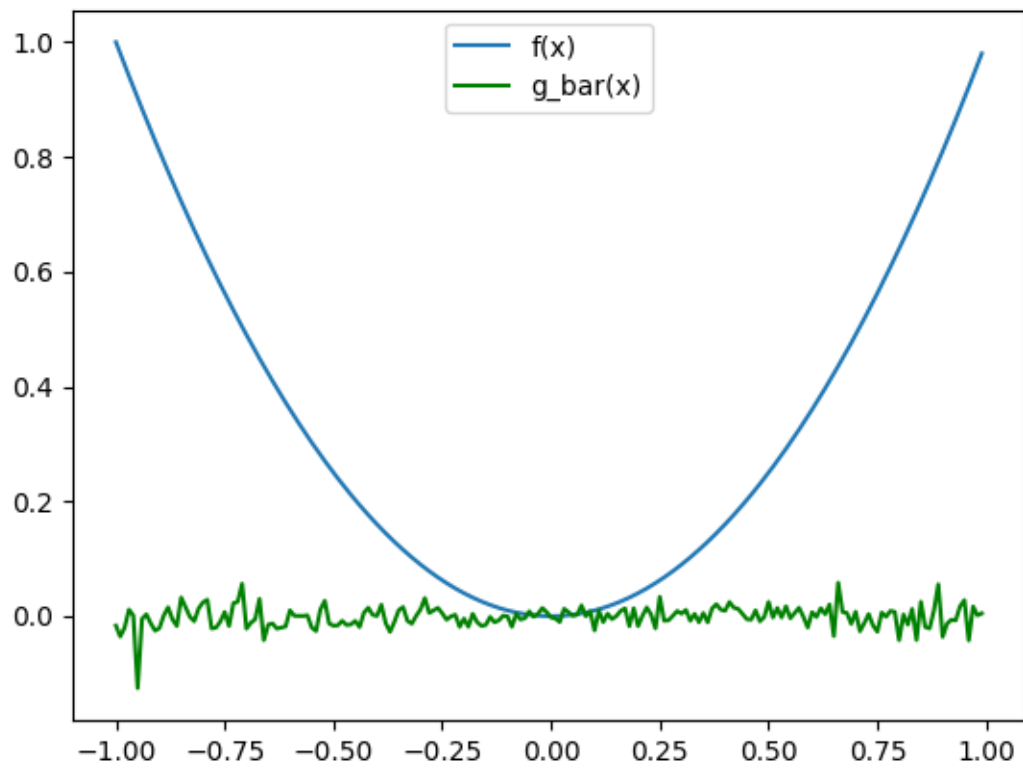
Learning Rate	$E_{in}$	Binary $E_{in}$	Binary $E_{out}$	Iterations	Time (s)
0.01	0.69	0.171	0.11	22429	0.49
0.1	0.663	0.171	0.11	2239	0.05
1	0.487	0.171	0.11	220	0.01
4	0.522	0.171	0.11	51	0.003
7	0.767	0.171	0.11	44	0.00299
7.5	0.813	0.171	0.11	180	0.00601
7.6	0.822	0.171	0.11	446	0.012
7.7	0.408	0.164	0.117	1000000	22.03

I normalized both the training and testing set to the mean and variance of the training set. This is because we never want to include the testing data in training.

Normalizing the data increases the performance of the model. Using the same parameters without normalized data yields  $E_{in}$  and  $E_{out}$  (binary) approaching random guessing.

Changing the learning rate usually arrives at the same error but just at different rates. The smallest learning rate takes > 20000 iterations to arrive at the same prediction model as does a learning rate of 7 (which takes just 44). Too big of a learning rate becomes cumbersome very quickly however. Increasing the learning rate to 7.6 makes us take 446 iterations and 7.7 makes us max out the number of iterations at  $10^6$ !

3 (c) plot.



2.

Problem 2.22

With noise  $E_{\text{out}}(g^{(D)}) = E_{x,y}[(g^{(D)}(\vec{x}) - y(\vec{x}))^2]$   
 where  $y(\vec{x}) = f(\vec{x}) + \epsilon$  ( $\epsilon$  zero mean w/var  $\sigma^2$ )

Show bias-variance decomp becomes

$$E_D[E_{\text{out}}(g^{(D)})] = \sigma^2 + \text{bias} + \text{var.}$$

- Redo book derivation with

$$\begin{aligned} E_D[E_{\text{out}}(g^{(D)})] &= E_D[E_{x,y}[(g^{(D)}(\vec{x}) - (f(\vec{x}) + \epsilon))^2]] \\ &= E_D[E_{x,y}[g^{(D)}(\vec{x})^2 - 2g^{(D)}(\vec{x})(f(\vec{x}) + \epsilon) + (f(\vec{x}) + \epsilon)^2]] \end{aligned}$$

- Switch  $E_D$  and  $E_{x,y}$ , distribute  $E_D$

$$= E_{x,y}[E_D[g^{(D)}(\vec{x})^2] - 2E_D[g^{(D)}(\vec{x})](f(\vec{x}) + \epsilon) + (f(\vec{x}) + \epsilon)^2]$$

- add in a  $-\bar{g}(\vec{x})^2 + \bar{g}(\vec{x})^2$

$$= E_{x,y}[E_D[g^{(D)}(\vec{x})^2] - \bar{g}(\vec{x})^2 + \bar{g}(\vec{x})^2 - 2\bar{g}(\vec{x})(f(\vec{x}) + \epsilon) + (f(\vec{x}) + \epsilon)^2]$$

$$= E_{x,y}[E_D[g^{(D)}(\vec{x})^2 - \bar{g}(\vec{x})^2] + \bar{g}(\vec{x})^2 - 2\bar{g}(\vec{x})f(\vec{x}) - 2\bar{g}(\vec{x})\epsilon + f(\vec{x})^2 + 2f(\vec{x})\epsilon + \epsilon^2]$$

$$= E_{x,y}[\text{variance} + \bar{g}(\vec{x})^2 - 2\bar{g}(\vec{x})f(\vec{x}) + f(\vec{x})^2 - 2\bar{g}(\vec{x})\epsilon + 2f(\vec{x})\epsilon + \epsilon^2]$$

$$= E_{x,y}[\text{variance} + (\bar{g}(\vec{x}) - f(\vec{x}))^2 + \epsilon^2 + (2f(\vec{x}) - \bar{g}(\vec{x}))\epsilon]$$

$$= \text{variance} + \text{bias} + E_{x,y}[\epsilon^2] - 2E_{x,y}[(\bar{g}(\vec{x}) - f(\vec{x}))\epsilon]$$

$$= \boxed{\text{variance} + \text{bias} + \sigma^2} \leftarrow$$

$\rightarrow$  this is 0 b/c  $E[\epsilon] = 0$

3. 2.24

$$x \in [-1, 1], f(x) = x^2$$

$$D = \{(x_1, x_1^2), (x_2, x_2^2)\}$$

$$h \in \{h(x) = ax + b\}$$

(a) Give an expression for  $\bar{g}(x)$  ok



$$a = \frac{x_2^2 - x_1^2}{x_2 - x_1} \quad b = \frac{x_2 x_1^2 - x_1 x_2^2}{x_2 - x_1}$$

Find  $E_D[ax + b]$

$$\bar{g}(x) = E_D \left[ \frac{x_2^2 - x_1^2}{x_2 - x_1} x + \frac{x_2 x_1^2 - x_1 x_2^2}{x_2 - x_1} \right]$$

$$= x \frac{1}{2} \frac{1}{2} \int_{-1}^1 \int_{-1}^1 \frac{x_2^2 - x_1^2}{x_2 - x_1} dx_1 dx_2 + \frac{1}{2} \frac{1}{2} \int_{-1}^1 \int_{-1}^1 \frac{x_1 x_2^2 - x_2 x_1^2}{x_2 - x_1} dx_1 dx_2$$

$$= \frac{x}{4} \int_{-1}^1 \int_{-1}^1 \frac{(x_2 - x_1)(x_2 + x_1)}{x_2 - x_1} dx_1 dx_2 + \frac{1}{4} \int_{-1}^1 \int_{-1}^1 \frac{-x_1 x_2 (x_2 - x_1)}{x_2 - x_1} dx_1 dx_2$$

$$= \frac{x}{4} \int_{-1}^1 \int_{-1}^1 (x_2 + x_1) dx_1 dx_2 + \frac{1}{4} \int_{-1}^1 \int_{-1}^1 -x_1 x_2 dx_1 dx_2$$

$$= \frac{x}{4} 0 + 0 + \frac{1}{4} 00 = \boxed{0}$$

(b) Numerically determine  $\bar{g}(x)$ ,  $E[E_{out}]$ , bias & Var

-  $\bar{g}(x)$ : generate  $N$  datasets  $((x_1, x_1^2), (x_2, x_2^2))$

• Calculate  $a = \frac{1}{N} \sum_{n=1}^N (x_{1n} + x_{2n})$

$$b = \frac{1}{N} \sum_{n=1}^N -x_{1n} x_{2n}$$

- bias, Var,  $E[E_{out}]$ :

• sample  $N$  datasets, calculate  $g^D(x)$  and  $\bar{g}(x)$

• resample above  $M$  times

• Variance:  $E[E_D[(g^D(x) - \bar{g}(x))^2]]$

$$= \frac{1}{M} \frac{1}{N} \sum_{m=1}^M \sum_{n=1}^N (g_m^D(x) - \bar{g}_m(x))^2$$



$$\bullet \text{Bias} = \mathbb{E}_x[\mathbb{E}_D[(\bar{g}(x) - f(x))^2]]$$

$$\text{Numerically: } \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N (\bar{g}_{m,n}(x) - f_{m,n}(x))^2$$

$$\bullet \mathbb{E}[E_{\text{out}}] = \mathbb{E}_x[\mathbb{E}_D[(g^D(x) - f(x))^2]]$$

$$\text{Numerically: } \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N (g_{m,n}^D(x) - f_{m,n}(x))^2$$

(c)

$$\text{Variance} = 0.336$$

$$\text{Bias} = 0.207$$

$$\mathbb{E}[E_{\text{out}}] = 0.544$$

pretty close

(figure on next page)

$$\text{Var} + \text{bias} = 0.543$$

(d)

Analytically compute  $\mathbb{E}[E_{\text{out}}]$ , bias, var.

$$\text{bias} = \mathbb{E}_x[\mathbb{E}_D[(\bar{g}(x) - f(x))^2]]$$

$$= \mathbb{E}_x[(0 - f(x))^2]$$

$$= \mathbb{E}_x[f(x)^2] = \mathbb{E}_x[x^4] = \frac{1}{2} \int_{-1}^1 x^4 dx = \frac{1}{2} \frac{x^5}{5} \Big|_{-1}^1$$

$$\text{bias} = \frac{1}{2} \frac{2}{5} = \frac{1}{5}$$

$$\text{Variance} = \mathbb{E}_x[\mathbb{E}_D[(g^D(x) - \bar{g}(x))^2]]$$

$$= \mathbb{E}_x[\mathbb{E}_D[(ax+b)^2]] = \mathbb{E}_x[\mathbb{E}_D[(x(x+x_2)x - x_1x_2)^2]]$$

$$= \mathbb{E}_x[\mathbb{E}_D[(x_1+x_2)^2 x^2 - 2x_1x_2(x_1+x_2)x + x_1^2x_2^2]]$$

$$= \mathbb{E}_x[x^2 \mathbb{E}_D[(x_1+x_2)^2] - 2x \mathbb{E}_D[x_1x_2(x_1+x_2)] + \mathbb{E}_D[x_1^2x_2^2]]$$

$$= \mathbb{E}_x[x^2 \mathbb{E}_D[x_1^2+x_2^2+2x_1x_2] - 2x \mathbb{E}_D[x_1x_2+x_2^2x_1] + \mathbb{E}_D[x_1^2x_2^2]]$$

$$\mathbb{E}_D[x_1] = 0, \mathbb{E}_D[x_1^2] = \frac{1}{2} \int_{-1}^1 x^2 dx = \frac{1}{2} \frac{x^3}{3} \Big|_{-1}^1 = \frac{1}{2} \frac{2}{3} = \frac{1}{3}$$

$$= \mathbb{E}_x[x^2(\frac{1}{3}+\frac{1}{3}) - 2x(\frac{1}{3}(0)+\frac{1}{3}(0)) + (\frac{1}{3}\frac{1}{3})]$$

$$= \frac{1}{3}(\frac{2}{3}) - 0 + \frac{1}{9} = \frac{3}{9} = \frac{1}{3} = \text{Variance}$$

(d) cont.

$$\begin{aligned} E[E_{\text{out}}] &= E_x[E_D[(g(x) - f(x))^2]] && \text{(forget about } E_D \text{ for a sec)} \\ &= E_x[(ax+b-x^2)^2] = E_x[x^4 - 2x^2(ax+b) + (ax+b)^2] \\ &= E_x[x^4] - 2aE_x[x^3] - 2bE_x[x^2] + a^2E[x^2] + 2abE[x] + b^2 \\ &= \frac{1}{5} - 0 - 2b\left(\frac{1}{3}\right) + a^2\left(\frac{1}{3}\right) + b^2 \\ &= \frac{1}{5} + \frac{1}{3}(a^2 - 2b) + b^2 \\ &= E_D\left[\frac{1}{5} + \frac{1}{3}(a^2 - 2b) + b^2\right] && a = x_1 + x_2, \quad b = -x_1x_2 \\ &= \frac{1}{5} + \frac{1}{3}E_D[(x_1+x_2)^2 + 2x_1x_2] + E_D[(x_1x_2)^2] \\ &= \frac{1}{5} + \frac{1}{3}E_D[x_1^2 + 2x_1x_2 + x_2^2 + 2x_1x_2] + E_D[x_1^2x_2^2] \\ &= \frac{1}{5} + \frac{1}{3}\left(\frac{1}{3} + 0 + \frac{1}{3} + 0\right) + \left(\frac{1}{3}\right)\frac{1}{3} \\ &= \frac{1}{5} + \frac{2}{9} + \frac{1}{9} = \boxed{8/15} = E[E_{\text{out}}] \quad (\text{also equals bias + var}) \end{aligned}$$



4. 3.4

# Adaptive Linear Neuron (adaline)

(a)  $e_n(\vec{w}) = (\max(0, 1 - y_n \vec{w}^T \vec{x}))^2$

• show  $e_n(\vec{w})$  is continuous and differentiable.

• Write down  $\nabla e_n(\vec{w})$

- 0 is continuous and differentiable ( $1 - y_n \vec{w}^T \vec{x} < 0$ )

-  $(1 - y_n \vec{w}^T \vec{x})^2$  is continuous and differentiable ( $1 - y_n \vec{w}^T \vec{x} > 0$ )

- if  $1 - y_n \vec{w}^T \vec{x} = 0$ , then  $y_n = \vec{w}^T \vec{x}$  and there is no error!

$\rightarrow e_n(\vec{w}) = 0$

$\rightarrow$  so this function is continuous between both its pieces

-  $\nabla e_n(\vec{w}) = 2(1 - y_n \vec{w}^T \vec{x})(-y_n \vec{x})$

- This is 0 at  $1 - y_n \vec{w}^T \vec{x}$ .

$\rightarrow$  so this function is differentiable between both its pieces.

(b) Show  $e_n(\vec{w})$  is an upper bound for  $\mathbb{I}[\text{sign}(\vec{w}^T \vec{x}_n) \neq y_n]$

Two cases:

1.  $\text{sign}(\vec{w}^T \vec{x}_n) = y_n$

- here  $y_n \vec{w}^T \vec{x} \geq 0$ ,  $\mathbb{I}[\text{sign}(\vec{w}^T \vec{x}_n) \neq y_n] = 0$ ,

-  $e_n(\vec{w}) \geq 0$  if  $y_n \vec{w}^T \vec{x} < 1$

-  $e_n(\vec{w}) = 0$  if  $y_n \vec{w}^T \vec{x} \geq 1$

- Thus:  $\mathbb{I}[\text{sign}(\vec{w}^T \vec{x}_n) \neq y_n] \leq e_n(\vec{w})$  (in this case)

2.  $\text{sign}(\vec{w}^T \vec{x}_n) \neq y_n$

- here  $y_n \vec{w}^T \vec{x} \leq 0$ ,  $\mathbb{I}[\text{sign}(\vec{w}^T \vec{x}_n) \neq y_n] = 1$

-  $e_n(\vec{w}) \geq 1$

- Thus  $\mathbb{I}[\text{sign}(\vec{w}^T \vec{x}_n) \neq y_n] \leq e_n(\vec{w})$  (in both cases now)

Thus  $E_n(\vec{w}) = \frac{1}{N} \sum_{n=1}^N \mathbb{I}[\text{sign}(\vec{w}^T \vec{x}_n) \neq y_n] \leq \frac{1}{N} \sum_{n=1}^N e_n(\vec{w})$

(c) Argue Adaline algorithm performs stochastic gradient descent on  $\frac{1}{N} \sum_{n=1}^N e_n(\vec{w})$ .

$$s(t) = \vec{w}^T(t) \vec{x}(t)$$

AD  $\vec{w}(t+1) \leftarrow \vec{w}(t) - \eta \nabla_{\vec{w}} e_n(\vec{w})$   $\rightarrow$  gradient descent

$$\vec{w}(t+1) \leftarrow \vec{w}(t) - \eta (2(1 - y_n \vec{w}^T \vec{x}) (-y_n \vec{x}))$$

$$\leftarrow \vec{w}(t) - 2\eta (1 - y_n s(t)) (-y_n \vec{x})$$

$$\leftarrow \vec{w}(t) + 2\eta (y_n - y_n s(t)) \vec{x}$$

$$\leftarrow \vec{w}(t) + 2\eta (y_n - s(t)) \vec{x}$$

$$\boxed{\vec{w}(t+1) \leftarrow \vec{w}(t) + 2\eta (y(t) - s(t)) \vec{x}(t)}$$

This is Adaline algorithm (substitute  $\eta = 2\eta$ ).



## 5.3.19 Is learning feasible?

$$(a) \quad \Phi(\vec{x}) = \begin{cases} (0, \dots, 0, 1, 0, \dots) & \text{if } \vec{x} = \vec{x}_n \\ (0, 0, \dots, 0) & \text{otherwise} \end{cases}$$

This maps points of dimension 2 to dimension  $N$  (where  $N$  is total number of points). This dramatically increases the generalization error as  $N$  increases.

$$E_{out}(g) \leq E_{in}(g) + O\left(\sqrt{d_{VC} \frac{\ln N}{N}}\right)$$

here  $d_{VC}$  will grow with  $N$ , so the generalization error increases with  $N$  ☹️

$$(b) \quad \phi_n(\vec{x}) = \exp\left(-\frac{\|\vec{x} - \vec{x}_n\|^2}{2\sigma^2}\right)$$

This seems ok. It seems to just normalize  $\vec{x}_n$  to  $\vec{x}$  on a standard  $2\sigma^2$ . Looks like the polar form of something. Since the dimensionality stays the same after the transform, the VC dimension should stay the same and generalization should decrease with  $N$ .

$$(c) \quad \phi_{ij}(\vec{x}) = \exp\left(-\frac{\|\vec{x} - (i,j)\|^2}{2\sigma^2}\right) \quad \begin{matrix} i \in 0, 100, \dots \\ j \in 0, 100, \dots \end{matrix}$$

This converts every datapoint of dimension 2 to dimension  $101 \times 101$ . Dramatically increasing the dimensionality of the dataset. This gives us the same problem as in (a).