

PS 2.

1. a)  $P(y; \lambda) = \frac{\lambda^y}{y!}$   
 $= \frac{1}{y!} \cdot \lambda^y$

For exponential form:  $P(y; \lambda) = b(y) e^{\eta^T T(y) - a(\eta)}$   
 $= b(y) \cdot \frac{e^{\eta^T T(y)}}{e^{a(\eta)}}$

$\therefore \underline{b(y) = \frac{1}{y!}}, \quad \underline{a(\eta) = \lambda = e^\eta}$

$\lambda^y = e^{\eta^T T(y)}$

$\ln \lambda^y = \eta^T T(y)$ .

$\eta^T T(y) = y \cdot \ln \lambda$

~~$T(y) = \ln \lambda \cdot (\eta^T)^{-1} y$~~

$T(y) = y$

$\eta = (\ln \lambda)^*$

$a(\eta) = \lambda \Rightarrow a(\eta) = e^\eta$

1.b) canonical response function

$g(\eta) = E[T(y); \eta]$

~~$E[T(y)] = \lambda$~~        $\eta = \ln \lambda$   
 $g(\eta) = \ln$ ,

$$1. (c) \quad \log P(\vec{y}^{(i)} | \vec{x}^{(i)}; \theta)$$

$$\ell(\theta) = \log P(\vec{y} | X; \theta)$$

$$= \log \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta)$$

$$= \log \prod_{i=1}^m \frac{1}{y^{(i)}!} \cdot e^{\ln \lambda \cdot y^{(i)} - e^{\eta}}$$

$$= \cancel{\log \prod_{i=1}^m \frac{1}{y^{(i)}!}}$$

$$= \sum_{i=1}^m \log \left( \frac{1}{y^{(i)}!} \cdot e^{\ln \lambda \cdot y^{(i)} - e^{\eta}} \right)$$

$$= \sum_{i=1}^m \log \frac{1}{y^{(i)}!} + \eta \cancel{\ln \lambda \cdot y^{(i)} - e^{\eta}}$$

$$= \sum_{i=1}^m \log \frac{1}{y^{(i)}!} + \theta^T x^{(i)} \cdot y^{(i)} - e^{\theta^T x^{(i)}}$$

$$\frac{\partial \ell(\theta)}{\partial \theta_j} = \sum_{i=1}^m (x_j^{(i)} \cdot y^{(i)} - e^{\theta^T x^{(i)}} \cdot x_j^{(i)})$$

$$= \sum_{i=1}^m (y^{(i)} - e^{\theta^T x^{(i)}}) \cdot x_j^{(i)}$$

$$\eta = \ln \lambda = \theta^T x \Rightarrow \lambda = e^{\theta^T x}$$

$$h(x) = E[\tau(y) | x] = \lambda = e^{\theta^T x} = g^{-1}$$

$$\therefore \frac{\partial \ell(\theta)}{\partial \theta_j} = \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

~~1. c) plot~~

1.(c)

$$\theta_j^{t+1} \leftarrow \theta_j^t + \eta \frac{\partial l(\theta)}{\partial \theta_j}$$

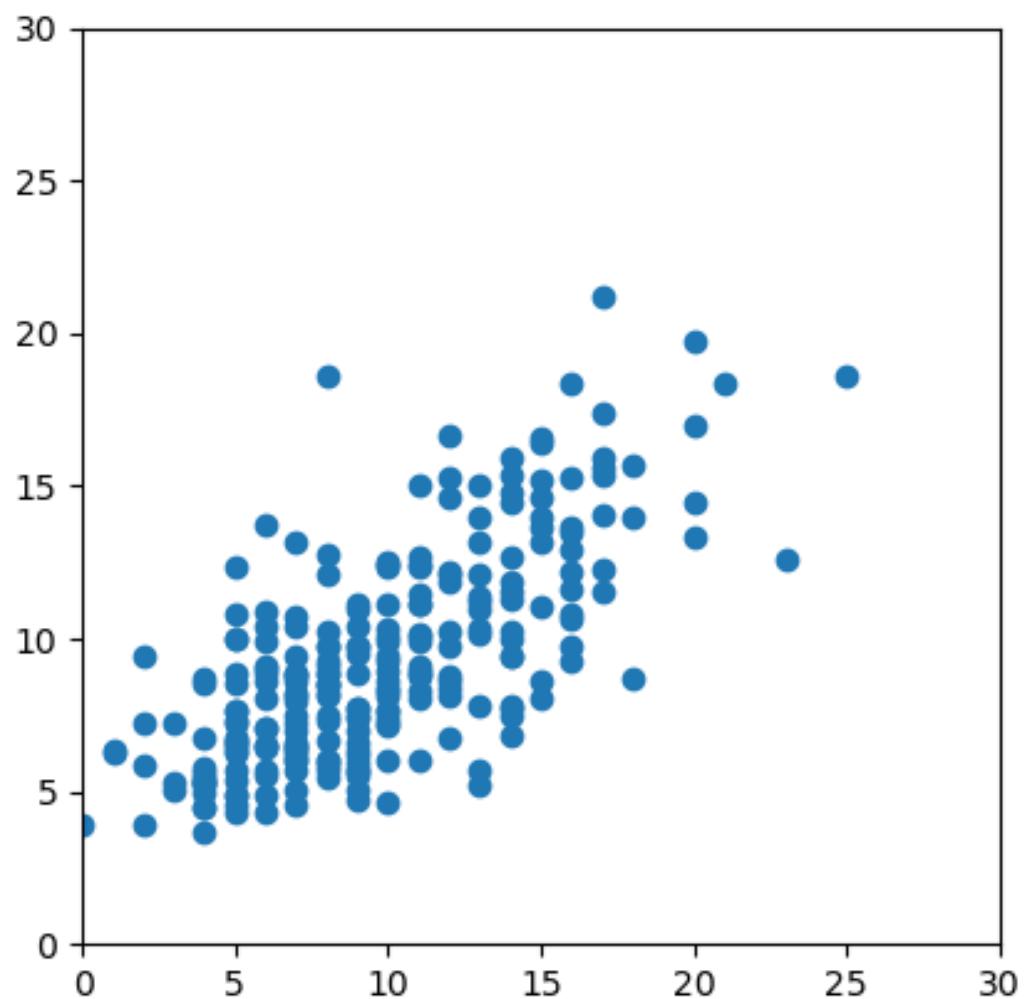
$$\theta^{t+1} \leftarrow \theta^t + \eta \frac{\partial l(\theta)}{\partial \theta_j}$$

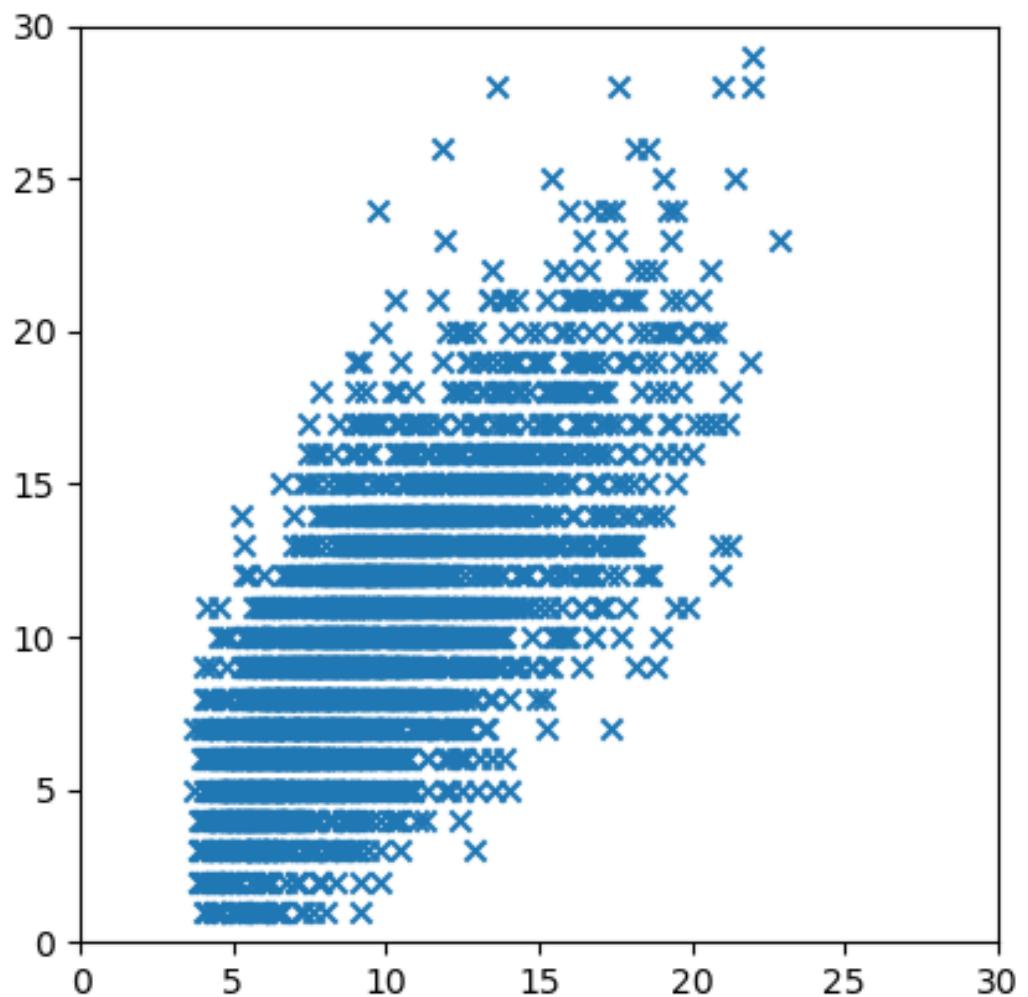
$$= \theta^t + \eta \cdot \sum_{i=1}^m (y^{(i)} - g^{-1}(x^{(i)})) \cdot x_j^{(i)}$$

1.(d) plot.

predict plot

validation plot





2. (a)

$$\int p(y; \eta) dy = 1 \quad \text{because the sum of all probability is 1}$$

$$\int b(y) \exp(\eta y - a(\eta)) dy = 1$$

$$\int \frac{b(y) \exp(\eta y)}{\exp(a(\eta))} dy = 1 \quad ; \because a(\eta) \text{ doesn't have } y.$$

$$\exp(a(\eta)) = \int b(y) \exp(\eta y) dy$$

$$a(\eta) = \log \int b(y) \exp(\eta y) dy$$

$$\begin{aligned}\frac{\partial a(\eta)}{\partial \eta} &= \frac{\partial}{\partial \eta} \log \int b(y) \exp(\eta y) dy \\ &= \frac{1}{\int b(y) \exp(\eta y) dy} \cdot \frac{\partial}{\partial \eta} \int b(y) \exp(\eta y) dy \\ &= \frac{1}{\exp(a(\eta))} \cdot \int \frac{\partial}{\partial \eta} b(y) \exp(\eta y) dy \\ &= \frac{1}{\exp(a(\eta))} \cdot \int b(y) \cdot \exp(\eta y) \cdot y dy \\ &= \int \frac{b(y) \exp(\eta y)}{\exp(a(\eta))} \cdot y dy \\ &= \int p(y; \eta) y dy \\ &= E(Y; \eta)\end{aligned}$$

2.(b)

$$\frac{\partial^2}{\partial \eta^2} a(\eta) = \frac{\partial}{\partial \eta} \left( \frac{\partial a(\eta)}{\partial \eta} \right)$$

$$= \frac{\partial}{\partial \eta} \int \frac{b(y) \exp(\eta y)}{\exp(a(\eta))} \cdot y \, dy$$

$$= \int b(y) y \cdot \frac{\partial \exp(\eta y - a(\eta))}{\partial \eta} \, dy$$

$$= \int b(y) y \cdot \exp(\eta y - a(\eta)) \cdot \left[ y - \frac{\partial a(\eta)}{\partial \eta} \right] \, dy$$

$$= \int b(y) \exp(\eta y - a(\eta)) y^2 \, dy$$

$$- \int b(y) \exp(\eta y - a(\eta)) y \cdot \frac{\partial a(\eta)}{\partial \eta} \, dy$$

$$= \bar{E}[y^2] - \frac{\partial a(\eta)}{\partial \eta} \cdot \int b(y) \exp(\eta y - a(\eta)) y \, dy$$

$$= \bar{E}[y^2] - \bar{E}[y] \cdot \bar{E}(y)$$

$$= \bar{E}[y^2] - \bar{E}[y]^2$$

$$= \text{Var}(Y; \eta)$$

2.(C) for a single example

$$\begin{aligned}
 NLL &= l(\theta) = -\log(L(\theta)) = -\log(P(y_i | x_i; \theta)) \\
 &= -\log(b(y_i) \exp(\eta_i y_i - a(\eta_i))) \\
 &= -[\log(b(y_i)) + \eta_i y_i - a(\eta_i)] \\
 &= a(\eta_i) - \eta_i y_i - \log(b(y_i))
 \end{aligned}$$

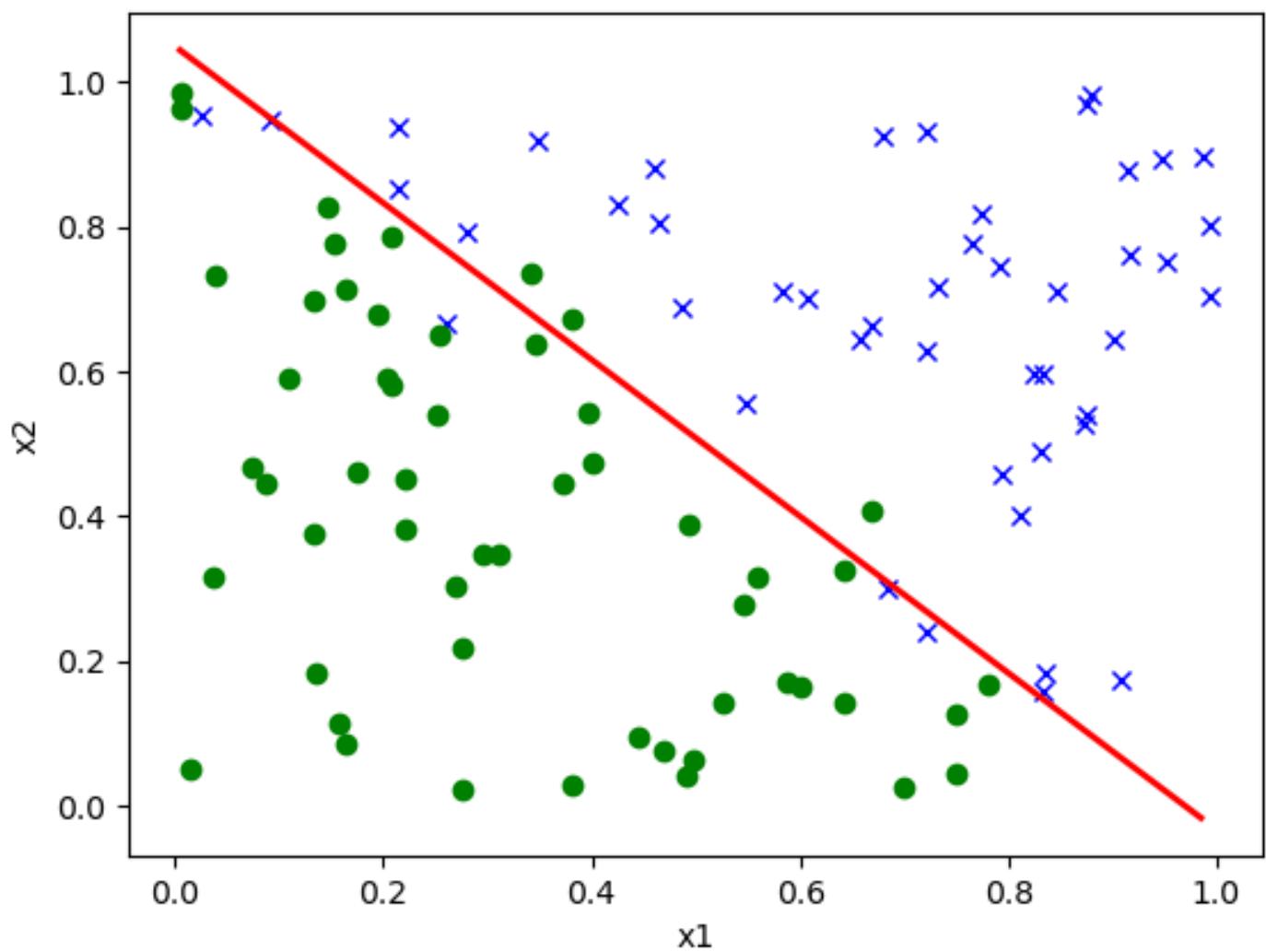
$$\text{Hessian} = \begin{bmatrix} \cdots & \cdots & \cdots \\ \vdots & \frac{\partial^2 l(\theta)}{\partial \theta_i \partial \theta_j} & \vdots \\ \cdots & \cdots & \cdots \end{bmatrix}$$

for a given  $i, j$ .

$$\begin{aligned}
 \frac{\partial^2 l(\theta)}{\partial \theta_i \partial \theta_j} &= \frac{\partial}{\partial \theta_i} \left( \frac{\partial l(\theta)}{\partial \theta_j} \right) = \frac{\partial}{\partial \theta_i} \left( \frac{\partial a(\eta_i)}{\partial \theta_j} - \frac{\partial \eta_i}{\partial \theta_j} y_i - \underbrace{\frac{\partial \log b(\eta_i)}{\partial \theta_j}}_{\downarrow 0} \right) \\
 &\stackrel{\eta_i = \theta^T x}{=} \frac{\partial}{\partial \theta_i} \left( x_i \frac{\partial a(\eta_i)}{\partial \eta_i} - x_i y_i \right) \left[ \begin{array}{l} \frac{\partial \eta_i}{\partial \theta_j} = x_j \\ \end{array} \right] \\
 &= \frac{\partial}{\partial \theta_i} \left( x_i \frac{\partial a(\eta_i)}{\partial \eta_i} \right) - \underbrace{\frac{\partial}{\partial \theta_i} x_i y_i}_{\rightarrow 0} \\
 &= x_i x_j \cdot \frac{\partial^2 a(\eta_i)}{\partial \eta_i^2} \rightarrow 0
 \end{aligned}$$

$$\text{Hessian} = \underbrace{\frac{\partial^2 a(\eta_i)}{\partial \eta_i^2}}_{\text{constant}} \cdot \begin{bmatrix} x_1^2 & x_1 x_2 & \dots & x_1 x_d \\ x_2 x_1 & x_2^2 & \dots & x_2 x_d \\ \vdots & & & \\ x_d x_1 & \dots & \dots & x_d^2 \end{bmatrix} = \frac{\partial^2 a(\eta_i)}{\partial \eta_i^2} \cdot \overline{X}^T \overline{X}$$

Where  $\overline{X} = [x_1, x_2, \dots, x_d]$  Hessian is by definition  $\bullet$  PSD =  $A^T A$  where  $A$  has independent col



3. (a) plot

(b) training set  $B$  ~~didn't~~ causes the training iteration maxed out  $10^5$  threshold: it didn't converge.  
its coefficient keeps increasing its magnitude.

because  $B$  is a linearly separable ~~not~~ data set.

it exists a  $\hat{w}$  such that  
for all positive training data

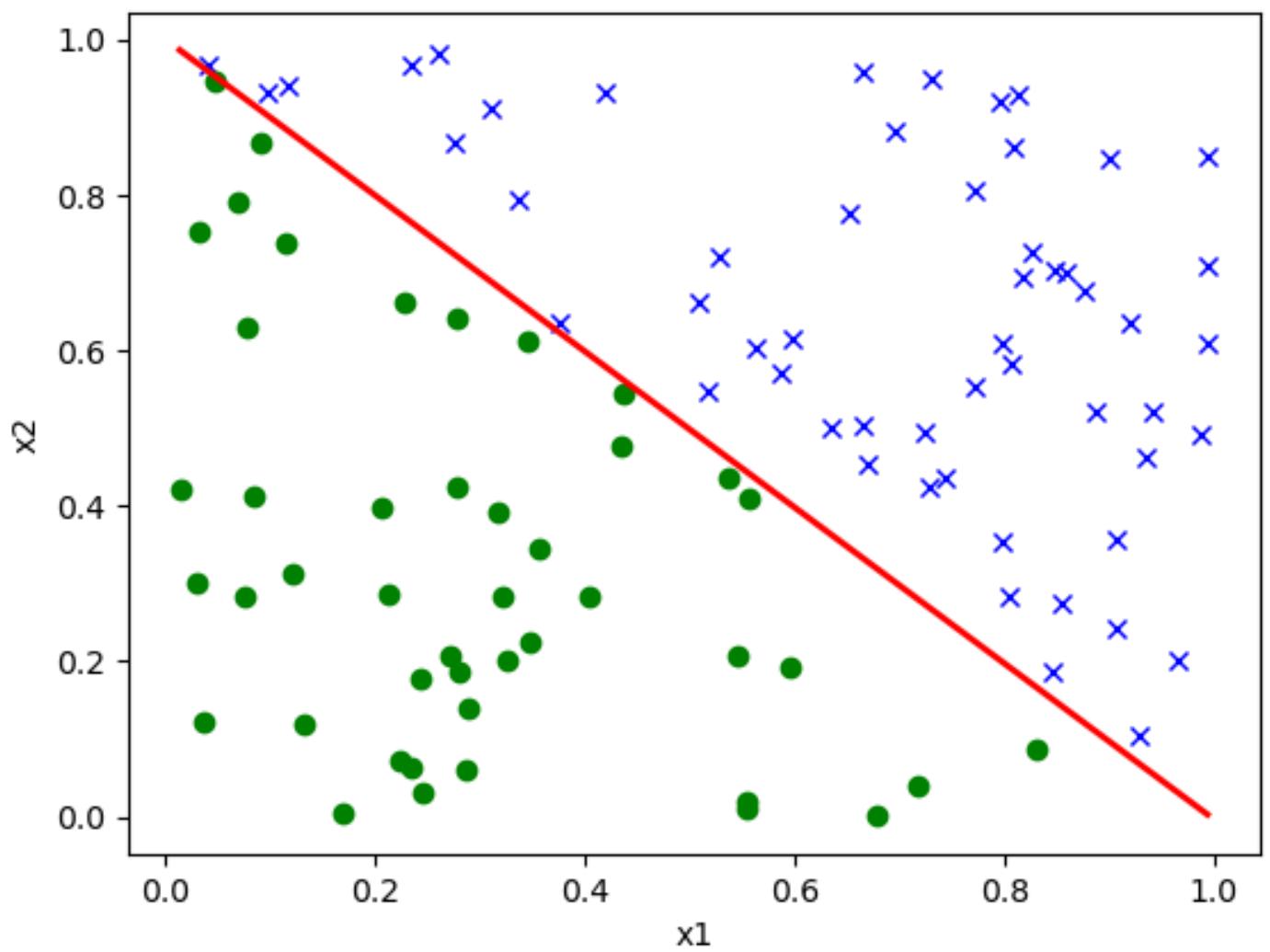
$$\hat{w}^T h(x) > 0$$

for all negative training data,  $\hat{w}^T h(x) < 0$ .

so,  $d\hat{w}$ , where  $d \in \mathbb{R}$  is also a decision boundary

because we maximize the likelihood estimation, we keep increasing the  $\hat{w}$  (overfitting), thus it keeps increasing.

A doesn't have this problem because it's not linearly separable in the ~~the~~ current dimension  $(X_1, X_2)$



3(c)-

- (i) no. it's linearly separable data.  
changing learning rate will not stop the coefficient from increasing.
- (ii) no. same as above
- (iii) no. linear scaling of linear separable data is still linear separable.
- (iv) yes,  $L_2$  regulation penalizes large coefficient so it will stop increasing
- v. yes. adding 0-mean Gaussian noise to data can change it to linear unseparable data.

3(d)

without Regularization:

$$\text{data set A: } \theta = [-20.7796 \quad 21.4171 \quad 2.3825]$$

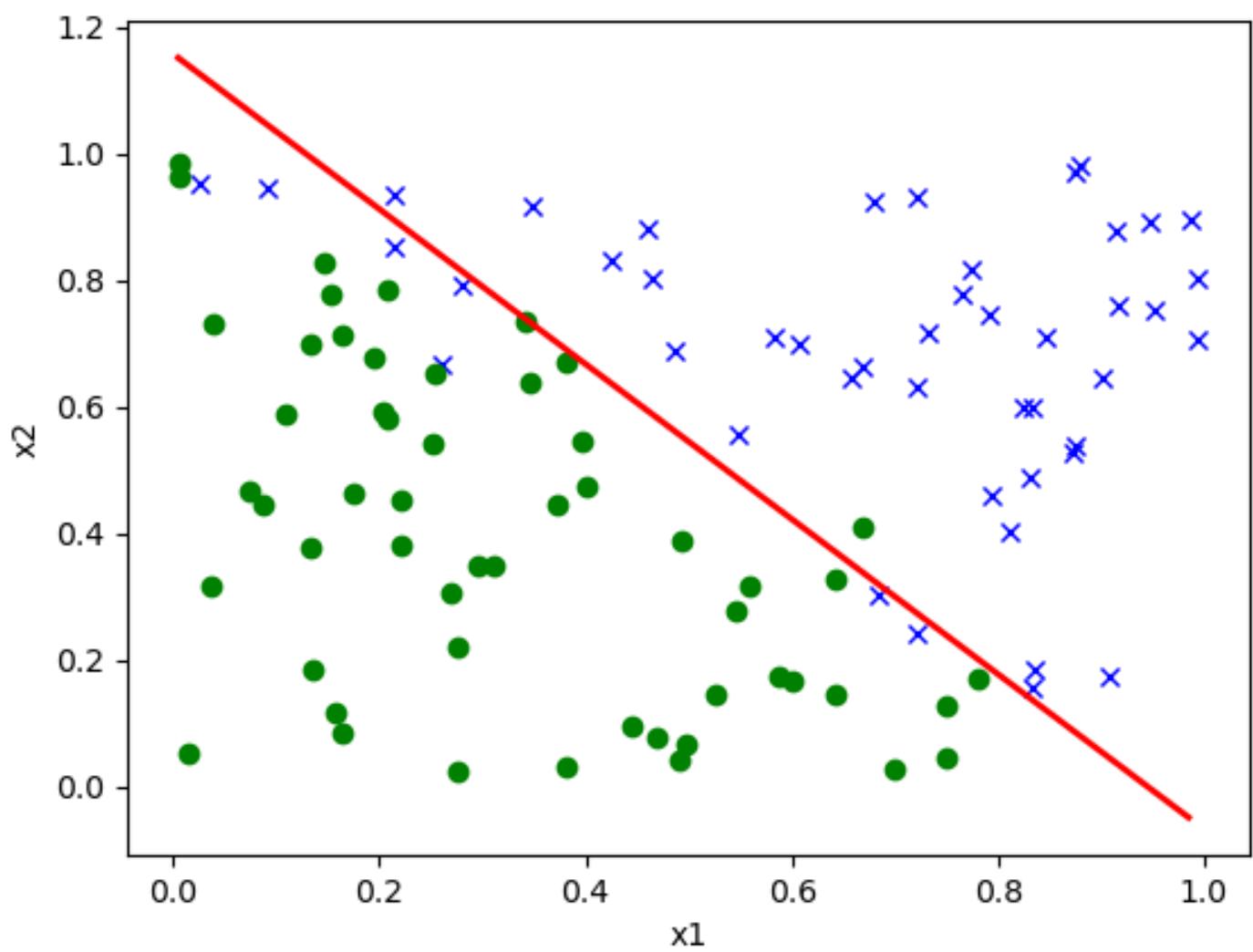
$$B: [-52.7200 \quad 52.9088 \quad 52.6758]$$

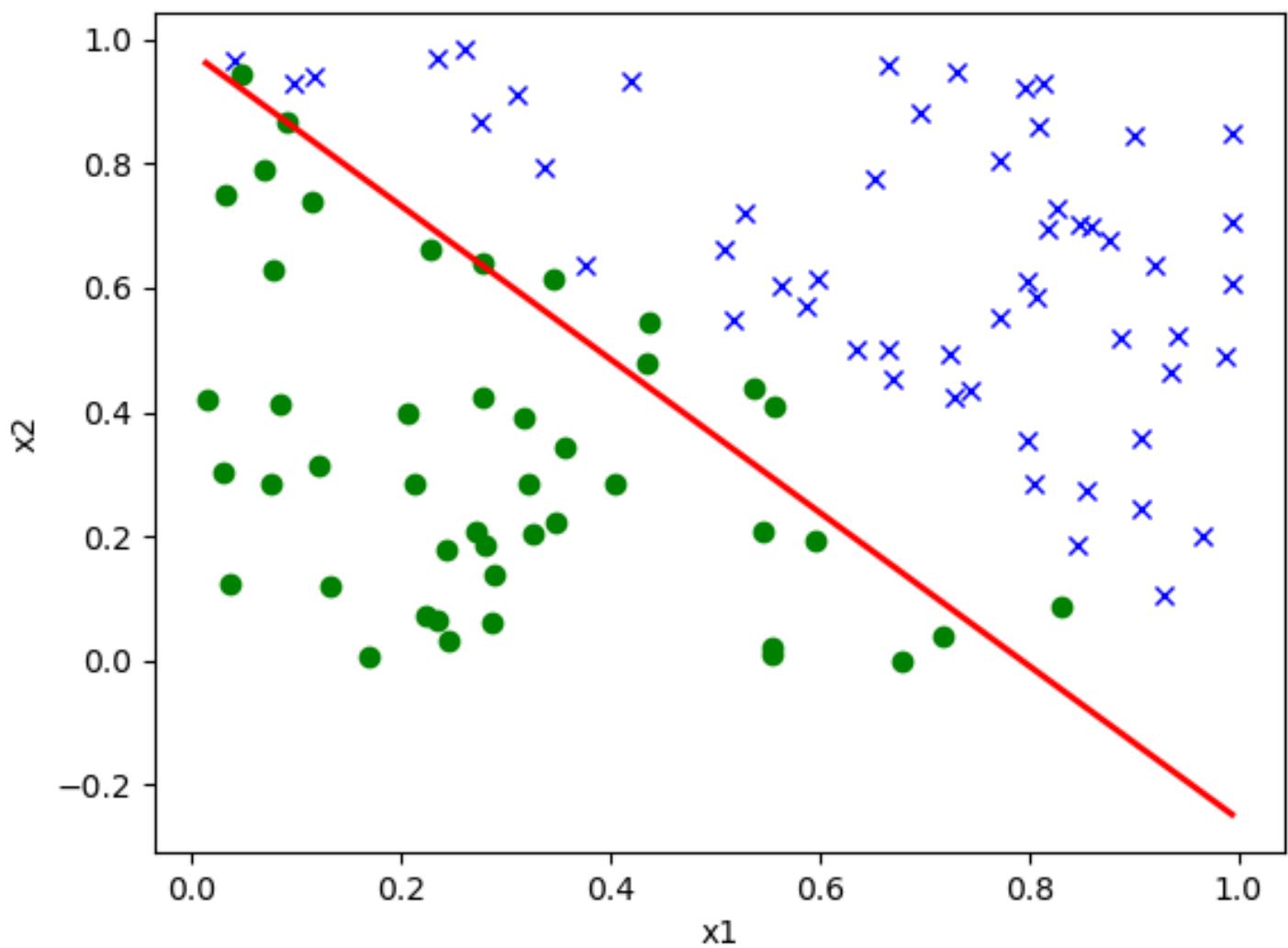
With Regularization.

$$A: [-2.5402 \quad 2.6892 \quad 2.1954]$$

$$B: [-2.3337 \quad 2.9438 \quad 2.3825]$$

plot





3.(d) (continue)

for A, regularization didn't significantly impact the decision boundary. However, for B, regularization moved the decision boundary to include some positive examples in negative prediction (i.e. false negative predictions).

Without regulation, the decision prediction will be overly confident. and for example, if such model is used to predict cancer, it might overfit to training data and gave false but confident prediction in real world. therefore, we'll prefer a regulated model to make sure it is not overfit.

4. (a) (i)

classifier = negative for all data.

then we will always be correct for (~~all~~) all the negative examples. ~~so~~  $1 - P$  is the accuracy.

$$\begin{aligned} \text{(ii)} \quad P &= \frac{P}{P+N} = \frac{TP+FN}{(TP+FN)+(TN+FP)} \\ &= \frac{TP+FN}{TP+TN+FP+FN} \end{aligned}$$

$$P_A + (1-P)A_0 = \frac{TP+FN}{TP+TN+FP+FN} \cdot \left( \frac{TP}{TP+FN} - \frac{TN}{TN+FP} \right) + \frac{TN}{TN+FP}$$

4.(a)(iii) continue.

$$\begin{aligned} & \rho A_1 + (1-\rho) A_0 \\ &= \frac{TP+FN}{S} A_1 + \frac{TN+FP}{S} A_0 \quad (\text{where } S = TP+TN+FP+FN) \\ &= \frac{TP+FN}{S} \cdot \frac{TP}{TP+FN} + \frac{TN+FP}{S} \cdot \frac{TN}{TN+FP} \\ &= \frac{TP}{S} + \frac{TN}{S} \\ &= \frac{TP+TN}{S} = \frac{TP+TN}{TP+TN+FP+FN} \\ &= A. \end{aligned}$$

4.(a)(iii)

$$\begin{aligned} \bar{A} &= \frac{1}{2} (A_0 + A_1) \\ &= \frac{1}{2} \left( \frac{TP}{TP+FN} + \frac{TN}{TN+FP} \right) \end{aligned}$$

for a classifier always predict false.

$$TP = 0$$

$$TN = (1-\rho) \cdot S$$

where  $S = TP+TN+FP+FN$

$$FP = 0$$

$$FN = \rho \cdot S$$

$$\begin{aligned} \bar{A} &= \frac{1}{2} \left( \frac{0}{0+S} + \frac{(1-\rho)S}{(1-\rho)S+0} \right) \\ &= \frac{1}{2} \quad 0 \quad 1 \end{aligned}$$

4.(b).

$$\bar{A} = 0.947273$$

$$\bar{A} = 0.7775$$

$$A_0 = 0.985$$

$$A_1 = 0.57$$

### Plot

4.(c), let  $\frac{\# \text{ of}}{\text{all examples in } D}$  be  $S$ . then  $N = (1-P)S$ ,  
 ~~$D = P = PS$~~  where  $N$  is # of all negative examples  
 ~~$N = N =$~~  and  $P$  is  $\sim \sim \sim$  positive  $\sim$

in  $D'$   $\Rightarrow TN' = TN, FP' = FP$

$$N' = N = (1-P)S$$

$$P' = \frac{1}{K} \cdot P = \frac{(1-P)}{K} \cdot PS = (1-P)S \Rightarrow TP' = \frac{1}{K} TP, FN' = \frac{1}{K} FN$$

$$\therefore N' = P'$$

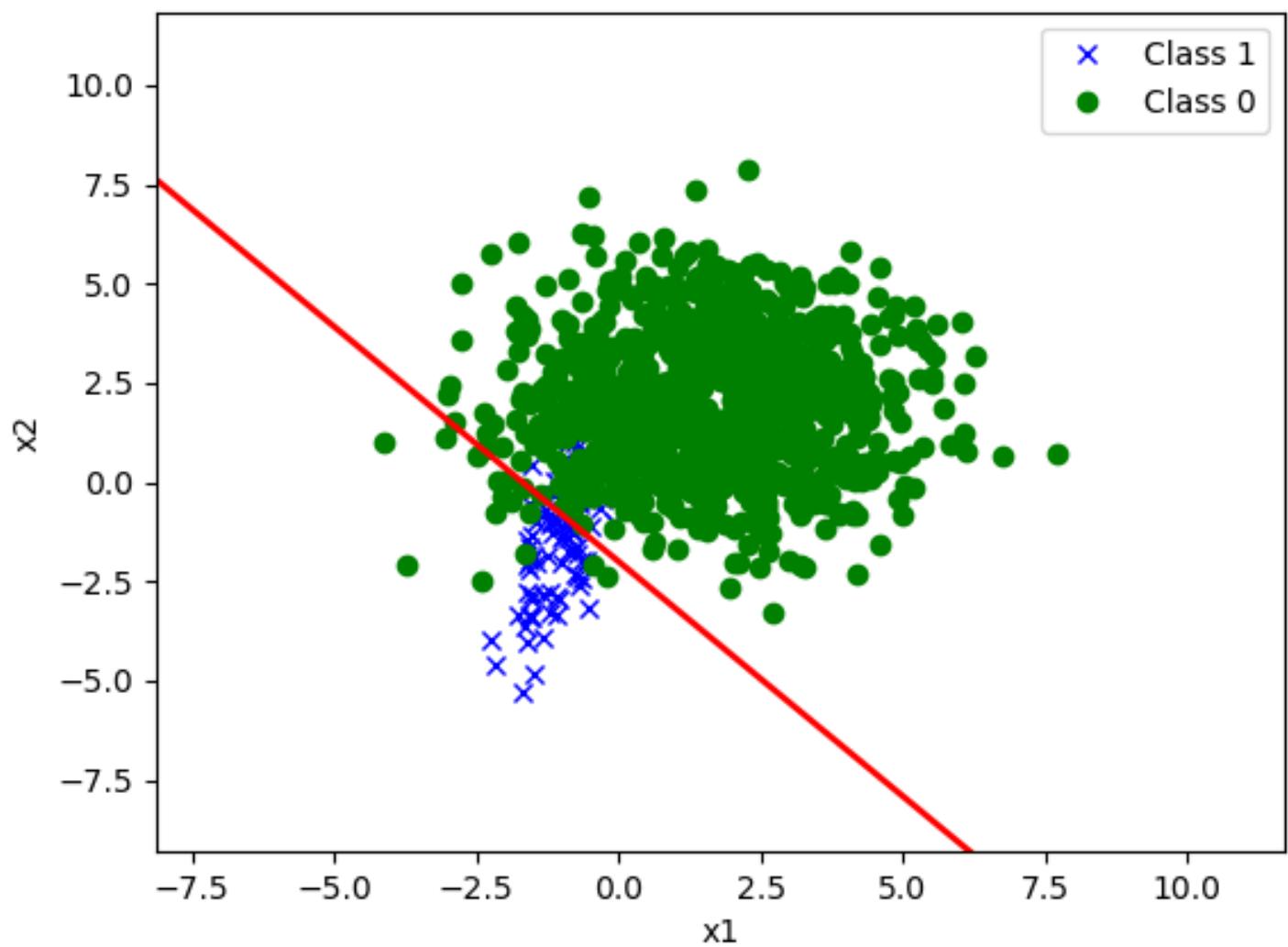
$$Q' = \frac{P'}{N'+P'} = \frac{1}{2}$$

$$A' (\text{accuracy of } D') = P' A'_1 + (1-P') A'_0$$

$$= \frac{1}{2} (A'_1 + A'_0)$$

$$= \frac{1}{2} \left( \frac{TP'}{P'} + \frac{TN'}{N'} \right)$$

$$= \frac{1}{2} \left( \frac{\frac{TP}{P}}{K} + \frac{TN}{N} \right) = \frac{1}{2} \left( \frac{TP}{P} + \frac{TN}{N} \right) = \bar{A}$$



4.(c) (continue)

# of example in  $D'$  is  $n'$

$$n' = P' + N' = 2(1-P)n \quad \text{where } n \text{ is # of examples in } D.$$

$$-\frac{1}{n'} = -\frac{1}{2(1-P)n} = -\frac{1+K}{2n}$$

$J(0)$  for  $D'$

$$= -\frac{1}{n'} \sum_{i=1}^{n'} (y^{(i)} \log(h_0(x^i)) + (1-y^{(i)}) \log(1-h_0(x^i)))$$

$$= -\frac{1}{n'} \left[ \underbrace{\sum_{\substack{i=P'_1, P'_2, \dots \\ \downarrow \\ \text{all positive examples in } D'}}^{P'} y^{(i)} \log(h_0(x^i))}_{\substack{\sum_{\substack{i=N'_1, N'_2, \dots \\ \downarrow \\ \text{all negative examples in } D'}}^{N'} (1-y^{(i)}) \log(1-h_0(x^i))}} \right]$$

$$= -\frac{1}{n'} \left[ \underbrace{\frac{1}{K} \sum_{\substack{i=P_1, P_2, \dots \\ \downarrow \\ \text{all positive examples in } D}}^{P} y^{(i)} \log(h_0(x^i))}_{\substack{\sum_{\substack{i=N_1, N_2, \dots \\ \downarrow \\ \text{all negative examples in } D}}^{N} (1-y^{(i)}) \log(1-h_0(x^i))}} \right]$$

$$= -\frac{1}{n'} \left[ \sum_{\substack{i=P_1, P_2, \dots}}^{P} \frac{1}{K} y^{(i)} \log(h_0(x^i)) + \sum_{\substack{i=N_1, N_2, \dots}}^{N} (1-y^{(i)}) \log(1-h_0(x^i)) \right]$$

$$= -\frac{1}{n'} \left[ \sum_{\substack{i=P_1, P_2, \dots}}^{P} w^{(i)} y^{(i)} \log(h_0(x^i)) + \sum_{\substack{i=N_1, N_2, \dots}}^{N} w^{(i)} (1-y^{(i)}) \log(1-h_0(x^i)) \right]$$

$w^{(i)} = 1$  if  $y^{(i)} = 0$ , i.e. negative examples.

$= \frac{1}{K}$  if  $y^{(i)} = 1$ ,  $\forall$  positive examples.

$$= -\frac{1}{n'} \sum_{i=1}^n w^{(i)} [y^{(i)} \log(h_0(x^i)) + (1-y^{(i)}) \log(1-h_0(x^i))]$$

$$= -\frac{1+K}{2n} \sum_{i=1}^n w^{(i)} [y^{(i)} \log(h_0(x^i)) + (1-y^{(i)}) \log(1-h_0(x^i))]$$

4.(d)

for D'

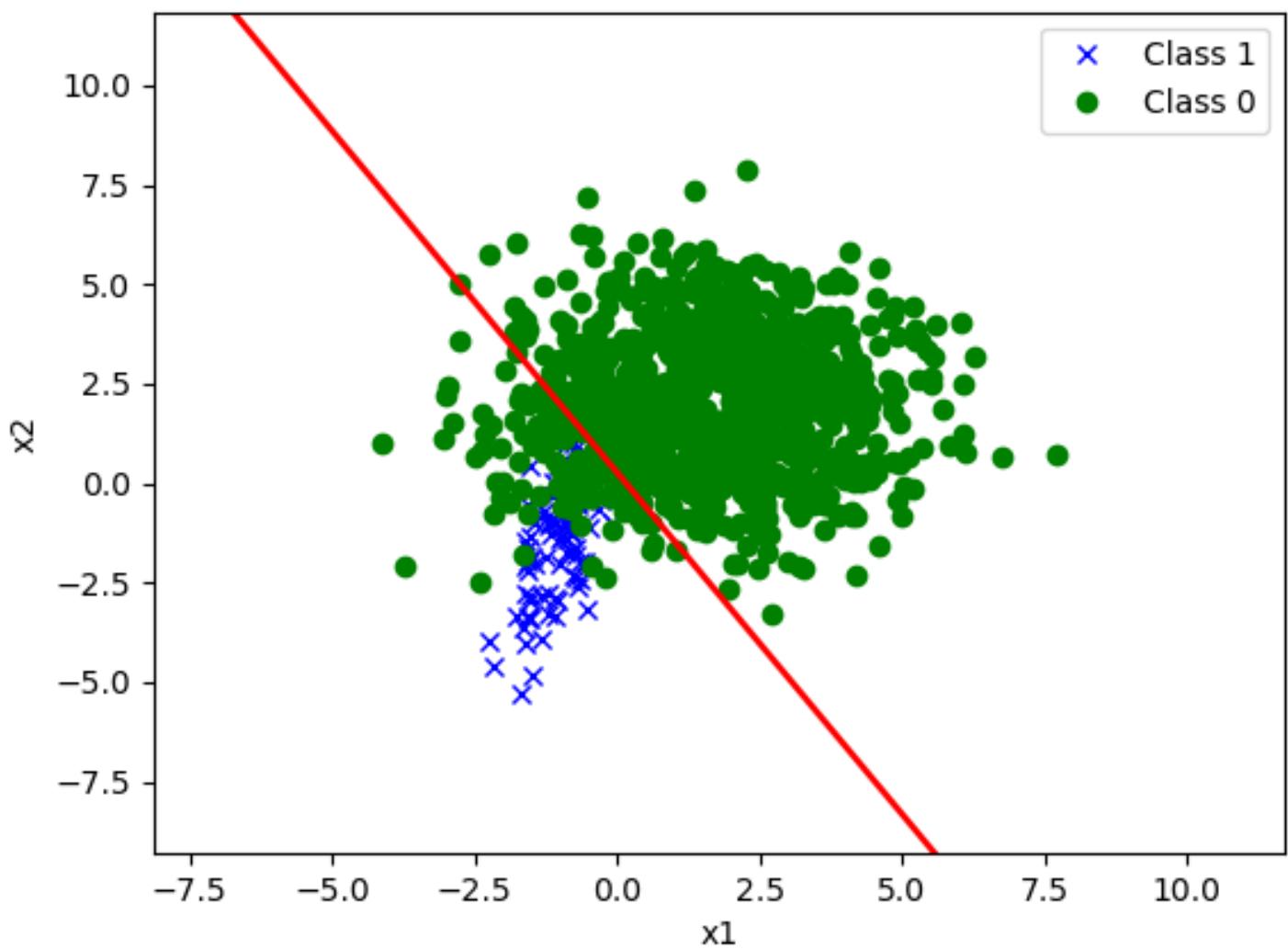
$$A = 0.8991$$

$$\bar{A} = 0.904$$

$$A_0 = 0.898$$

$$A_1 = 0.91$$

Plot



(5) (a)

$$\frac{\partial l_{CE}(t, y)}{\partial t} = \frac{\partial}{\partial t} \left[ -\log \left( \frac{\exp(t_y)}{\sum_{s=1}^K \exp(t_s)} \right) \right]$$

$$= \frac{\partial}{\partial t} \left[ -(\log(\exp(t_y)) - \log(\sum_{s=1}^K \exp(t_s))) \right]$$

$$= \frac{\partial}{\partial t} \left[ \log(\sum_{s=1}^K \exp(t_s)) - t_y \right]$$

$$= \begin{bmatrix} \frac{\partial}{\partial t_1} \left\{ \log \sum_{s=1}^K \exp(t_s) \right\} - t_y \\ \vdots \\ \frac{\partial}{\partial t_K} \left\{ \log \sum_{s=1}^K \exp(t_s) \right\} - t_y \end{bmatrix} . \text{ by chain rule:}$$

$$= \begin{bmatrix} \frac{\frac{\partial}{\partial t_1} \left( \sum_{s=1}^K \exp(t_s) \right)}{\sum_{s=1}^K \exp(t_s)} - \frac{\partial}{\partial t_1} t_y \\ \vdots \\ \frac{\frac{\partial}{\partial t_K} \left( \sum_{s=1}^K \exp(t_s) \right)}{\sum_{s=1}^K \exp(t_s)} - \frac{\partial}{\partial t_K} t_y \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\exp(t_1)}{\sum_{s=1}^K \exp(t_s)} - 0 \\ \frac{\exp(t_y)}{\sum_{s=1}^K \exp(t_s)} - 1 \end{bmatrix}$$

$$= \left[ \begin{array}{c} \frac{\exp(t_1)}{\sum_{s=1}^k \exp(t_s)} - 0 \\ \vdots \\ \frac{\exp(t_y)}{\sum_{s=1}^k \exp(t_s)} - 1 \\ \vdots \\ \frac{\exp(t_k)}{\sum_{s=1}^k \exp(t_s)} - 0 \end{array} \right] \leftarrow y^{\text{th}} \text{ item}$$

by definition.

$$= \text{softmax}(t) - e_y$$

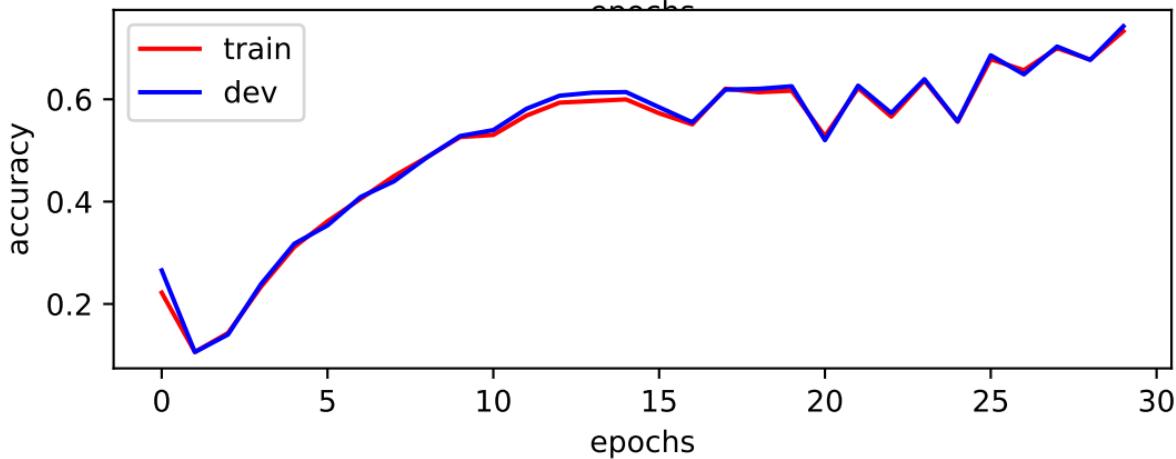
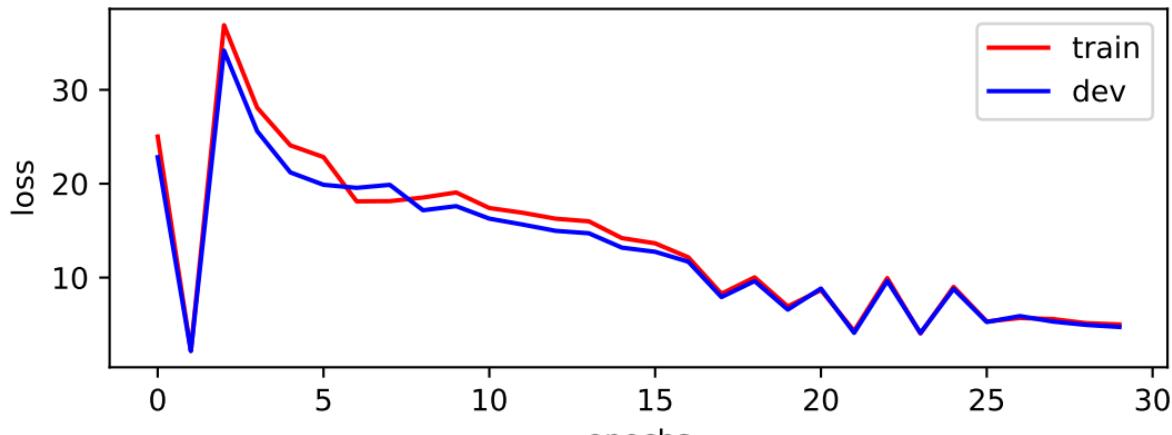
(5) b, plot

(c) plot

For model baseline, got accuracy: 0.733700

For model regularized, got accuracy: 0.633900

## Without Regularization



## With Regularization

