

* QUESTION 1

** (a) ① $E[L(Y=\text{keep}, t)]$

$$= \sum_{t \in T} L(Y=\text{keep}, t) \cdot \Pr(Y=\text{keep}, t)$$

$$= 0 \cdot 0.9 + 1 \cdot 0.1$$

$$= 0.1$$

They're
ind

② $E[L(Y=\text{remove}, t)]$

$$= \sum_{t \in T} L(Y=\text{remove}, t) \cdot \Pr(Y=\text{remove}, t)$$

$$= 100 \cdot 0.9 + 0 \cdot 0.1$$

$$= 90$$

** (b) \vec{X} : Feature Vector

Since we know $\Pr(t=\text{spam}|\vec{X})$, we can calculate $\Pr(t=\text{non-spam}|\vec{X}) = 1 - \Pr(t=\text{spam}|\vec{X})$.

* Bayes' Optimal Decision Rule:

Find a y^* from {Keep, Remove}, that minimizes

week 2 p.9.52 $E[L(Y=y^*, t)]$, as follows:

$$\left\{ \begin{array}{l} y^* = \text{keep}: E[L(Y=\text{keep}, t)] = 1 \cdot \Pr(t=\text{spam}|\vec{X}) + 0 \cdot [1 - \Pr(t=\text{spam}|\vec{X})] \\ y^* = \text{remove}: E[L(Y=\text{remove}, t)] = 0 \cdot \Pr(t=\text{spam}|\vec{X}) + 100 \cdot [1 - \Pr(t=\text{spam}|\vec{X})] \end{array} \right.$$

Plug in $\Pr(t=\text{spam}|\vec{X})$, compare the two values, choose the one that outputs lower expectation.

** (c) First, let's compute: # Note $\Pr(t=\text{spam}) = 0.1$

$\Pr(t=\text{spam}, X_1=x_1, X_2=x_2)$	$X_2=0$	$X_2=1$
$X_1=0$	$\frac{200}{4691}$	$\frac{103}{103}$
$X_1=1$	$\frac{200}{209}$	1

$$\Pr(X_1=0, X_2=0 | t=\text{spam}) = 0.4$$

$$\Pr(X_1=0, X_2=0 | t=\text{non-spam}) = 0.998$$

$$\Rightarrow \Pr(t=\text{spam} | X_1=0, X_2=0) = \frac{\Pr(X_1=0, X_2=0 | t=\text{spam}) \cdot \Pr(t=\text{spam})}{\sum_{t \in T} \Pr(X_1=0, X_2=0 | t) \cdot \Pr(t)}$$

Similarly,

$$\Pr(t=\text{spam} | X_1=0, X_2=1) = \frac{103}{209} \quad \# \text{ As in Table 1, so we know } \Pr(t=\text{spam}|\vec{X}), \text{ we also can get}$$

$$\Pr(t=\text{spam} | X_1=1, X_2=0) = \frac{200}{209}$$

so, we can find expected Loss of Keep and Remove as follows: $\Pr(t=\text{non-spam}|\vec{X})$, AND then we can apply optimal rule from (b).

$X_2=0$	$X_2=1$
0.0426	0.4709
95.7365	2.9126

$X_1=0$	$X_1=1$
0.4569	1
4.3062	0

Hence, we have the following actions based on x_1 and x_2 .

If $\begin{cases} x_1=0, x_2=0 \rightarrow \text{keep } (0.0426) \\ x_1=0, x_2=1 \rightarrow \text{keep } (0.9709) \\ x_1=1, x_2=0 \rightarrow \text{keep } (0.9569) \\ x_1=1, x_2=1 \rightarrow \text{remove } (0) \end{cases}$

***(d) Based on (c), we know four best actions for four possible combinations of x_1 and x_2 .

$$\text{So, } E[L(Y_k, t)] = 0.0426 + 0.9709 + 0.9569 + 0 \\ = 1.9704$$

$$\frac{\partial L}{\partial w_j} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_j}$$

$$\frac{\partial L}{\partial z} = \left\{ \begin{array}{ll} *11 & , 20, 11 \end{array} \right.$$



*QUESTION 2

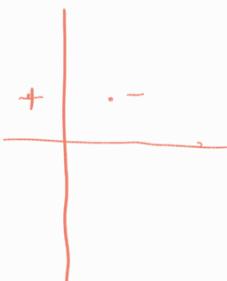
**(a) proof:



- If two points lie in a half-space, the line segment connecting them also lie in the same half-space.
- suppose some feasible weights exist ...
- but $-1 \rightarrow 3$ passes $1, 1$ exists in another half-space \rightarrow contradiction.

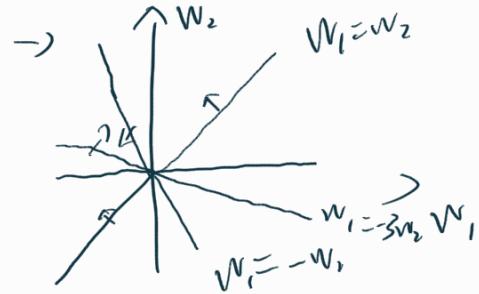
**(b) The feature maps:

x	$\psi_1(x)$	$\psi_2(x)$	t
-1	-1	1	1
1	1	1	0
3	3	9	1



(Hard threshold)

$$\Rightarrow \begin{cases} -w_1 + w_2 \geq 0 \\ w_1 + w_2 < 0 \\ 3w_1 + 9w_2 \geq 0 \end{cases}$$



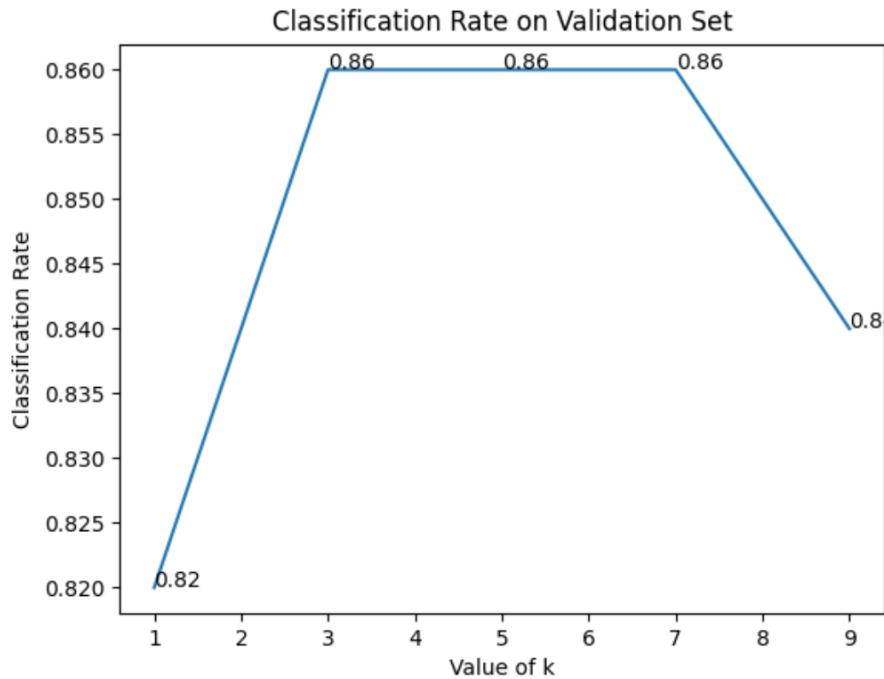
AND a pair $(-1, \frac{1}{2})$ will classify all the examples correctly.

$$\begin{aligned} -(-1) + \frac{1}{2} &= \frac{3}{2} \geq 0 \\ -1 + \frac{1}{2} &= -\frac{1}{2} < 0 \\ 3 \times (-1) + 9 \times \frac{1}{2} &= \frac{3}{2} \geq 0 \end{aligned}$$

* Question 3

** Part(1)

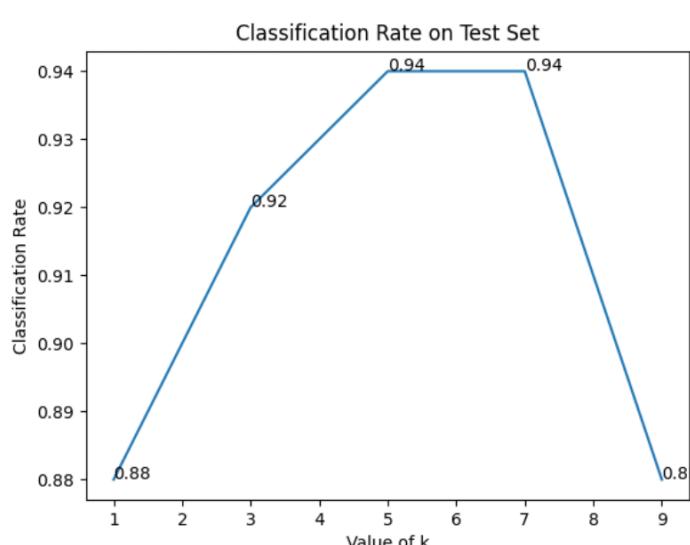
*** (a)



*** b)

OVERALL, the classifier performs good on validation set the $k \in \{1, 3, 5, 7, 9\}$. As shown in the plot, the classifier predicts correctly for around 82% data points.

However, different k s do perform differently. I would like to choose $k^* = 5$, the classification rate for k^* is 0.86. AND the classification rate for $k^* + 2 (7)$ is 0.86, $k^* - 2 (3)$ is 0.86.



The test performance of these values of K are similar to that from validation:

- ① $K=5, 7$ are still the best choices.
- ② $K=3$ is not the best choice anymore, but still a better choice than $K=1$ and $K=9$.

Such findings are expected, because we know when $K \rightarrow 0$, the model underfits, and when $K \rightarrow \infty$, the model overfits, both won't result in best rate, a middle value of K is therefore preferred.

**Part(Z)

*** (a) Implemented and uploaded to markus.

*** (b) Run `check_grad` is very small as follows:

```
[[ -0.10461503 -0.10461503]
 [-0.2001114 -0.20011139]
 [-0.20279696 -0.20279696]
 [-0.19851759 -0.19851759]
 [-0.13922366 -0.13922366]
 [ 0.05024501  0.05024502]
 [ 0.27444364  0.27444364]
 [ 0.12492774  0.12492773]
 [ 0.02181367  0.02181367]
 [-0.05183629 -0.05183628]
 [ 0.01158533  0.01158533]]
diff = 2.6438378464677305e-08
```

I'd assume 2.64×10^{-8} as
"small enough".

THE BEST HYPERPARAMETERS for two datasets.

- `load_train()`:

learning-rate: 0.3

num-iterations: 200

CE on train: 0.1789

classification rate on train (errtot): 0.9875 (0.015)

CE on validation: 0.2928

classification rate on validation (errtot): 0.9 (0.1)

- `load_train_small()`:

learning-rate: 0.05

num-iterations: 50

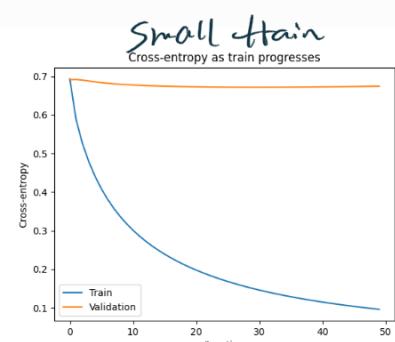
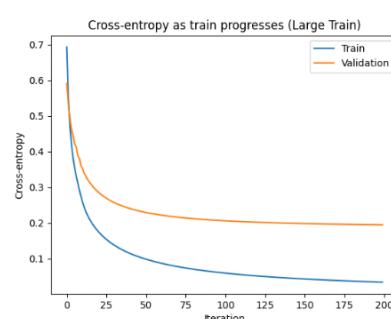
CE on train: 0.0941

classification rate on train (errtot): 1

CE on validation: 0.6743

classification rate on validation (errtot): 0.64 (0.36)

*** (c) Look at the following two images.



* I run my code several times with diff hyperparameters, the results do change a little bit, and I'll select based on the following rules. { ① better on validation, to avoid overfit.
② the curve that changes faster → to avoid long time of tuning.

*QUESTION 4

** (a) I'll apply "Direct Solution":

I'll use ξ to denote cost function.

So, we can compute:

$$\xi = \frac{1}{2} \sum_{i=1}^N a^i (y^i - \sum_{j=1}^N w^j \cdot x_j^i)^2 + \frac{\lambda}{2} \sum_{j=1}^N w_j^2$$

$$\begin{aligned} \frac{\partial \xi}{\partial w_k} &= \frac{1}{2} \sum_{i=1}^N a^i (y^i - \sum_{j=1}^N w^j \cdot x_j^i) \cdot (-x_k^i) + \frac{\lambda}{2} \cdot 2 \cdot w^k \\ &= \sum_{i=1}^N a^i \cdot x_k^i (\sum_{j=1}^N w^j \cdot x_j^i - y^i) + \lambda \cdot w^k \end{aligned}$$

We want $\frac{\partial \xi}{\partial w_k} = 0$, so:

$$\lambda \cdot w^k = \sum_{i=1}^N a^i \cdot x_k^i (y^i - \sum_{j=1}^N w^j \cdot x_j^i)$$

$$\lambda \cdot w^k = A \cdot x_k \cdot (Y - W^* \cdot X^T)$$

#Assume $X = \begin{bmatrix} \dots \\ \vdots \\ x_k \\ \vdots \\ \dots \end{bmatrix}$, $Y = \begin{bmatrix} \dots \\ \vdots \\ y \\ \vdots \\ \dots \end{bmatrix}$

$$\begin{aligned} \Leftrightarrow \lambda \cdot I \cdot w^* &= A \cdot X^T (Y - W^* \cdot X^T) \quad \# X^T \text{ because we want } A \text{ apply to } x_i \text{ in one sample.} \\ \Leftrightarrow \lambda \cdot I \cdot w^* &= A \cdot X^T Y - A \cdot X^T \cdot X \cdot w^* \end{aligned}$$

$$\Leftrightarrow \lambda \cdot I \cdot w^* + A \cdot X^T \cdot X \cdot w^* = A \cdot X^T Y$$

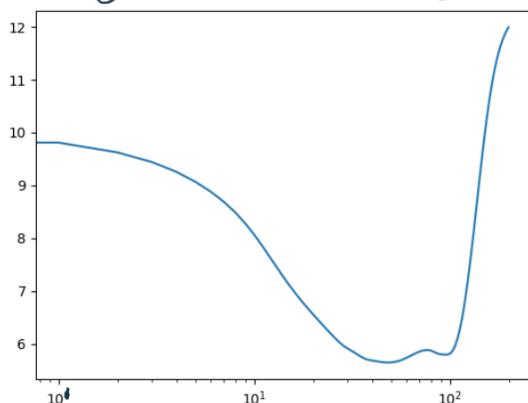
$$\Leftrightarrow \lambda \cdot I \cdot w^* + X^T A \cdot X \cdot w^* = X^T \cdot A \cdot Y \quad \# \text{ since } A \text{ is scalar matrix and } A_{ii} = a^i > 0,$$

$$\Leftrightarrow W^* = (\lambda I + X^T A X)^{-1} \cdot X^T A Y \quad \text{so } A X^T = X^T A$$

$$\Leftrightarrow W^* = (X^T A X + \lambda I)^{-1} \cdot X^T A Y$$

** (b)

Average Loss as Tau grows



** (c) When $\tau \rightarrow \infty$, the algorithm will perform bad, resulting in big loss in average. When $\tau \rightarrow 0$, it will predict badly and have a big loss as well.

It happened actually. For example, when

$\tau = 10$ and $\tau = 1000$, the average is biggest, and

there's a trend $\tau = 10 \xrightarrow{\text{decrease}} \text{a middle point} \xrightarrow{\text{increase}} \text{loss}$ around 50