

Q₁ a)

$$E[L(y=\text{keep}, t)] = 0.9 \cdot 0 + 0.1 \cdot 1 = 0.1$$

$$E[L(y=\text{remove}, t)] = 0.9 \cdot 100 + 0.1 \cdot 0 = 90;$$

Q₁b) suppose $P(t=\text{spam} | x)$ is known; then $P(t=\text{non-spam} | x) = 1 - P(t=\text{spam} | x)$

$$\text{then } E[L(y=\text{keep}, t)] = 1 \cdot P(t=\text{spam} | x) + 0 \cdot P(t=\text{non-spam} | x)$$

$$E[L(y=\text{remove}, t)] = 0 \cdot P(t=\text{spam} | x) + 100 \cdot P(t=\text{non-spam} | x)$$

thus $y^* = \begin{cases} \text{keep, if } E[L(y=\text{keep}, t)] < E[L(y=\text{remove}, t)] \\ \text{remove, otherwise} \end{cases}$

Q₁c) there are 4 pairs of feature vector: $(x_1, x_2) \rightarrow (0, 0), (0, 1), (1, 0), (1, 1)$

Since y^* considers $E[L(y, t)]$ and take the action $y \in \{\text{keep, remove}\}$ with minimal loss,

then:

$(0, 0)$: first infer $P(t|x)$ from $P(x|t)$

since $P(t|x) = \frac{P(x|t)P(t)}{P(x)}$, will find $P(x)$

$$P(x) = \sum_t P(x|t)P(t) = P(0, 0|\text{spam}) \cdot P(\text{spam}) + P(0, 0|\text{nonspam}) \cdot P(\text{nonspam})$$

$$= 0.4 \cdot 0.1 + 0.998 \cdot 0.9 = 0.9382$$

$$\text{thus } P(\text{spam} | (0, 0)) = \frac{0.4 \cdot 0.1}{0.9382} = 0.0426$$

$$P(\text{nonspam} | (0, 0)) = \frac{0.998 \cdot 0.9}{0.9382} = 0.957;$$

$$\text{then: } E[L(\text{keep}, t)] = 1 \cdot 0.0426 \quad \text{thus by } y^* \text{ from part b,}$$

$$E[L(\text{remove}, t)] = 100 \cdot 0.957 = 95.7 \quad y = \text{keep};$$

$$(0,1): p(x) = 0.3 \cdot 0.1 + 0.001 \cdot 0.9 = 0.0309$$

$$p(\text{spam} | (0,1)) = \frac{0.3 \cdot 0.1}{0.0309} = 0.971 \quad p(\text{nonspam} | (0,1)) = \frac{0.001 \cdot 0.9}{0.0309} = 0.029$$

$$E[L(\text{keep}, t)] = 1 \cdot 0.971 \quad E[L(\text{remove}, t)] = 100 \cdot 0.029 = 2.9$$

thus by y^* , $y = \text{keep}$

$$(1,0): p(x) = 0.2 \cdot 0.1 + 0.001 \cdot 0.9 = 0.0209$$

$$p(\text{spam} | (1,0)) = \frac{0.02}{0.0209} = 0.957 \quad p(\text{nonspam} | (1,0)) = \frac{0.001}{0.0209} = 0.043$$

$$E[L(\text{keep}, t)] = 0.957 \quad E[L(\text{remove}, t)] = 4.3$$

thus by y^* , $y = \text{keep}$;

$$(1,1): p(x) = 0.1 \cdot 0.1 + 0 \cdot 0.9 = 0.01$$

$$p(\text{spam} | (1,1)) = 1 \quad p(\text{nonspam} | (1,1)) = 0$$

$$E[L(\text{keep}, t)] = 1 \quad E[L(\text{remove}, t)] = 0$$

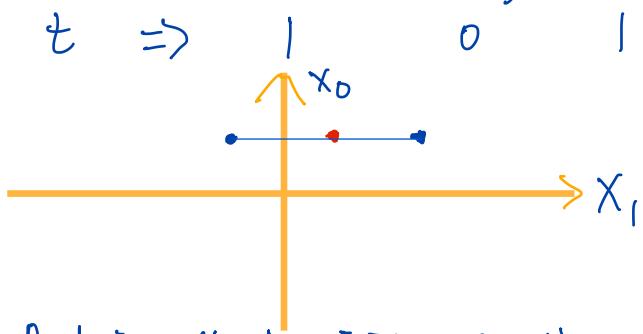
thus by y^* , $y = \text{remove}$

$$Q(d) > E[L(y^*, t)] = 1 \cdot \underset{\text{keep}}{p((0,0) | \text{spam})} \cdot p(\text{spam}) + 1 \cdot p((0,1) | \text{spam}) \cdot$$

$$p(\text{spam}) + 1 \cdot p((1,0) | \text{spam}) \cdot p(\text{spam}) + 100 \cdot p((1,1) | \text{nonspam}) \cdot p(\text{nonspam})$$

$$= 1 \cdot 0.04 + 1 \cdot 0.03 + 1 \cdot 0.02 + 100 \cdot 0 = 0.09$$

Q2 a) $(x_0, x_1) \Rightarrow (1, -1), (1, 1), (1, 3)$



Realizing the line joining $(1, -1), (1, 3)$ also contains $(1, 1)$;
thus the point $(1, 1)$ is supposed to lie in a different
region but result in the same area as $(1, -1)$ and $(1, 3)$.
Thus not linearly separable.

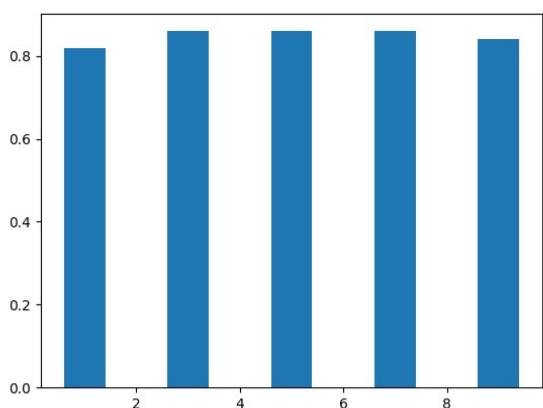
Q2 b) after applying feature vector, will acquire below table:

$$\begin{array}{ccc} x & x^2 & t \\ -1 & 1 & 1 \\ 1 & 1 & 0 \\ 3 & 9 & 1 \end{array} \Rightarrow \begin{array}{l} -w_1 + w_2 \geq 0 \\ w_1 + w_2 < 0 \\ 3w_1 + 9w_2 \geq 0 \end{array}$$

thus one choice of (w_1, w_2) is: $(-1, 0.5)$

Q3 KNN:

classification accuracy

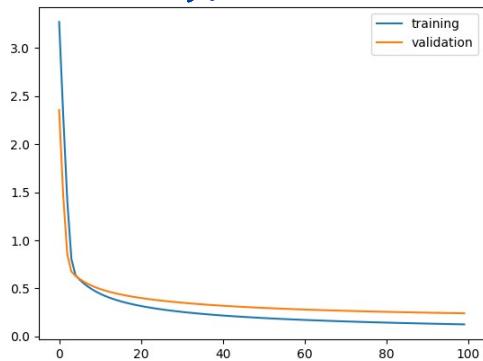


| bs KNN has highest accuracy when $k=5$;
| thus $k=5$ is preferred. The rate of accuracy
| is around 0.86;
| when $k=3$ and $k=7$, the accuracy is also 0.86;
| after running knn with $k=5$ on test set,
| the accuracy is 0.94, which is higher
| than validation set's performance.

Q3 logistic regression:

C>

mnist 'ce'



$$\lambda = 0.1$$

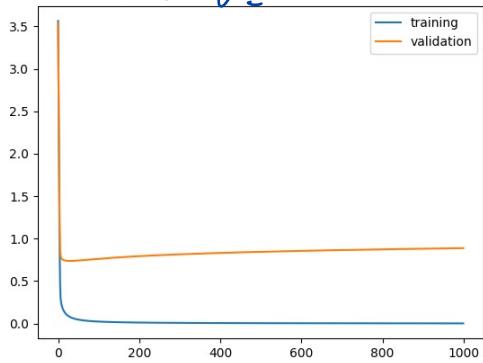
$$\text{num_iter} = 100$$

final training ce: 0.1260

final validation ce: 0.2411

final test ce: 0.2312

mnist_small_ce



$$\lambda = 0.1$$

$$\text{num_iter} = 1000$$

final training ce: 0.002

final validation ce: 0.8879

final test ce: 0.8672

several rounds of re-running doesn't result in significant changes to 'ce' tendency and final 'ce'.

Q4 a) $X: N \boxed{D}$

$W: D \boxed{1}$

Take derivative with respect to w_p ; since $(a-b)^2 = (b-a)^2$

$$\frac{\partial J}{\partial w_p} = \frac{\partial}{\partial w_p} \left[\frac{1}{2} \sum_{i=1}^N \sum_{d=1}^D a^i (-y^i + \sum_{d=1}^D w_d x_d^i)^2 + \frac{\lambda}{2} \sum_{d=1}^D (w_d)^2 \right]$$

$$= \left[\sum_{i=1}^N a^i (y^i + w_1 x_1^i + w_2 x_2^i + \dots + w_p x_p^i + \dots + w_D x_D^i) (x_p^i) \right] + \lambda w_p$$

$$= \frac{\textcircled{1}}{\textcircled{2}} W_1 \left(\sum_{i=1}^N a^i x_1^i x_p^i \right) + \dots + w_p \left(\sum_{i=1}^N a^i x_p^i x_p^i \right) + \dots + w_D \left(\sum_{i=1}^N a^i x_D^i x_p^i \right)$$

$$+ \lambda w_p - \frac{\sum_{i=1}^N a^i y^i x_p^i}{\textcircled{3}}$$

one observation: $a^i x_p^i = 0 \cdot x_p^1 + 0 \cdot x_p^2 + \dots + a^i \cdot x_p^i + \dots + 0 \cdot x_p^n$
 $= [0 \dots a^i \dots 0] \cdot [x_p^1 \dots x_p^i \dots x_p^n]^T$

each $\frac{\partial J}{\partial w_p}$ is a sum of:

① product of: pth row of X^T matrix, ith row of diagonal matrix A, matrix X, and W^T vector.

② : pth row of λI multiply with W^T

③ ith row of diagonal matrix A with vector y and pth column of matrix X;

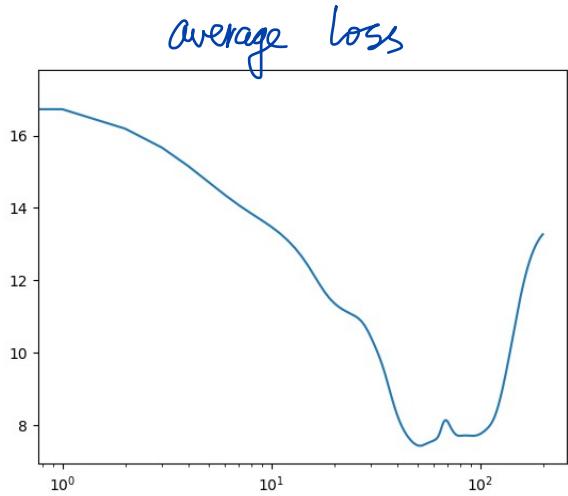
Thus $\frac{\partial J}{\partial W} = \begin{bmatrix} \vdots \\ \frac{\partial J}{\partial w_p} \\ \vdots \end{bmatrix} = -(A^T y)^T + X^T (AX) \cdot W + \lambda I W = 0$

\Rightarrow since A is diagonal, then $(Ay)^T = y^T A^T = y^T A$

$$(X^T AX + \lambda I) W = (y^T A X)^T = X^T A y$$

thus $W = (X^T A X + \lambda I)^{-1} X^T A y$

Q4 c>



as $T \rightarrow \infty$, the loss of this algorithm would be higher than 7, the minimal loss ; same for the case when $T \rightarrow 0$.

by running simulation with $T = 2, 3, 1000000, 10^8$,

the result is: when $T > 10^5$, the loss stabilizes at 13.5 above when T is 2, 3, the loss is 25 above, which are both higher than 7, the minimal loss could achieve.