

ESE 542 FINAL EXAMProblem 1

$$P(Y=0) = P$$

$$P(Y=1) = 1 - P$$

$$P(X=x | Y=0) = \frac{1}{S} \text{ for } -4 \leq x < 1$$

$$P(X=x | Y=1) = \frac{1}{4} \text{ for } 0 \leq x \leq 4$$

$$P(X=x | Y=y)$$

a)

10

$$0.25 = \frac{1}{4}$$

$$P(Y=x | Y=1)$$

$$P(Y=x | Y=0)$$

$$\frac{1}{S} = 0.2$$

15

-4

1

4

x

b)

By Bayes' rule,

20

$$P(Y=y | X=x) = \frac{P(X=x | Y=y) \cdot P(Y=y)}{P(X=x)}$$

25

30

$$c) h^*(x) = \underset{y \in \{0, 1\}}{\operatorname{argmax}} P(Y=y | X=x)$$

$$= \underset{y \in \{0, 1\}}{\operatorname{argmax}} \frac{P(X=x | Y=y) \cdot P(Y=y)}{P(X=x)} \quad (\text{Bayes' Rule})$$

$$= \underset{y \in \{0, 1\}}{\operatorname{argmax}} P(X=x | Y=y) \cdot P(Y=y) \quad (\text{Since } P(Y=y) \text{ doesn't depend on } y)$$

$$= \underset{y \in \{0, 1\}}{\operatorname{argmax}} \left\{ P(X=x | Y=0) \cdot P(Y=0), P(X=x | Y=1) \cdot P(Y=1) \right\}$$

$$= \underset{y \in \{0, 1\}}{\operatorname{argmax}} \left\{ \frac{1}{5} \cdot P, \frac{1}{4} \cdot (1-P) \right\}$$

$$= \begin{cases} 0 & \text{if } \frac{P}{5} > \frac{1-P}{4} \\ 1 & \text{otherwise} \end{cases}$$

Let's look @ the condition for 0:

$$\frac{P}{5} > \frac{1-P}{4} \quad \therefore h^*(x) = \begin{cases} 0 & \text{if } P > 5/9 \\ 1 & \text{if } P \leq 5/9 \end{cases}$$

$$\Rightarrow 4P > 5(1-P)$$

$$\Rightarrow 4P > 5 - 5P$$

$$\Rightarrow 9P > 5$$

$$\Rightarrow P > 5/9$$

This makes intuitive sense because we simply return the argument for y that maximizes $P(X=x | Y=y)$.

Note $5/9 > 1/2$, which means that if on average $Y=0$ is more likely, always output 0 irrespective of x .

$$a) P(h^*(x) \neq y)$$

$$= P(h^*(x) = 0, y=1) + P(h^*(x) = 1, y=0)$$

$$= P(h^*(x) = 0 | y=1) \cdot P(y=1) + P(h^*(x) = 1 | y=0) \cdot P(y=0)$$

$$= P(p > s_{1/2}) \times 1-p + P(p \leq s_{1/2}) \cdot p$$

$$= \mathbb{1}(p > s_{1/2}) \cdot (1-p) + \mathbb{1}(p \leq s_{1/2}) \cdot p$$

Therefore, the error rate is either p or $1-p$, depending on whether p is less than or greater than $s_{1/2}$ respectively.

This makes intuitive sense, because if we always classify 0, then we make an error whenever the label is 1, and that happens with $P(Y=1) = 1-p$.

On the other hand, if we always classify as 1, then we make an error whenever the label is 0, and that happens w/ ~~p~~ $P(Y=0) = p$.

In order to activate only one of p or $1-p$ as errors, they are multiplied with the indicator fns of their corresponding classifier conditions.

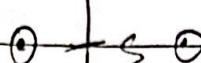
Legend

Problem 2

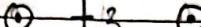
\square = Updated cluster center, O = Data point

\diamond = Optimal cluster center, \triangle = Initial cluster center

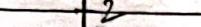
a)



5



10



15



20



25



b)

Points closest to $(-5, 5) \rightarrow \{(-5, 1), (-3, 1), (-1, 1), (1, 1)\}$

|| || || Points closest to $(-5, -5) \rightarrow \{(-5, -1), (-3, -1), (-1, -1), (1, -1)\}$

|| || || Points closest to $(5, -5) \rightarrow \{(5, -1), (3, -1), (1, -1), (-1, -1)\}$

|| || || Points closest to $(5, 5) \rightarrow \{(5, 1), (3, 1), (1, 1), (-1, 1)\}$

$$\text{Computing new center for first cluster} = \left(\frac{-5 + -3 + -1 + 1}{4}, \frac{1 + 1 + 3 + 5}{4} \right) \\ = \left(\frac{-10}{4}, \frac{10}{4} \right) = (-2.5, 2.5)$$

Similarly New center for second cluster = $(-2.5, -2.5)$

|| || || Third || = $(2.5, -2.5)$

|| || || fourth || = $(2.5, 2.5)$

c) Final output of clusters remains the same as those after termination

This is because after the first iteration, the new sets of clusters (ie points associated with each cluster center) remain the same.

As a result, the new cluster centers are the centroids of the same set of points from the previous iteration, and the new cluster centers don't change.

$$\therefore C_1 = C_2 = \text{Termination} = \{ (-2.5, 2.5), (-2.5, -2.5), (2.5, -2.5), (2.5, 2.5) \}$$

} Marked on plot with 

d) The optimal cluster centers should be

$$C_{OPT} = \{ (-4, 0), (0, -4), (4, 0), (0, 4) \}$$

} Marked on plot with 

Therefore in this instance, k-means does not find the optimal clusters

Problem 3

$$S = \{(x_i, y_i)\}_{i=1}^n$$

5) $x_i \in \mathbb{R}^d$

$$y_i = \beta^T x_i + \epsilon_i$$

$$\epsilon_i \in \mathbb{R}$$

$$\beta = (\beta_1, \dots, \beta_d)$$

14

This notation

is used

a) $\epsilon_i \stackrel{iid}{\sim} \text{Laplace}(0, 1) = \frac{1}{2} e^{-|z|} = l(z)$ later

Show: $f(y_1 - y_n | \beta) = \left(\frac{1}{2}\right)^n \prod_{i=1}^n e^{-|y_i - \beta^T x_i|}$

Now, we know that n

$$f(y_1 - y_n | \beta) = \prod_{i=1}^n f(y_i | \beta) \quad (\text{Since } \epsilon_i \text{ iid features also iid})$$

In class (cf. in Piazza @ 809), we saw that

$$\begin{aligned} \text{Now, } f(y_i | \beta) &= P(Y_i = y_i | \beta) \\ &= P(y_i = \beta^T x_i + \epsilon_i | \beta) \\ &= P(\epsilon_i = y_i - \beta^T x_i | \beta) \\ &= l(y_i - \beta^T x_i) \quad (\text{Since } \epsilon_i \text{ is Laplace}(0, 1) \text{ distributed}). \end{aligned}$$

$$\therefore f(y_1 - y_n | \beta) = \prod_{i=1}^n f(y_i | \beta)$$

$$= \prod_{i=1}^n l(y_i - \beta^T x_i)$$

$$\Rightarrow \prod_{i=1}^n \left(\frac{1}{2}\right) e^{-|y_i - \beta^T x_i|} \quad \begin{array}{l} \text{Since } \frac{1}{2} \text{ does} \\ \text{not depend on} \\ \text{index } i \end{array}$$

$$= \left(\frac{1}{2}\right)^n \prod_{i=1}^n e^{-|y_i - \beta^T x_i|}$$

b) We have now established that

$$f(y_i - y_n | \beta) = L(\beta) = \left(\frac{1}{2}\right)^n \prod_{i=1}^n e^{-|y_i - \beta^T x_i|}$$

$$\therefore \hat{\beta}_{MLE} = \operatorname{argmax}_{\beta} f(y_i - y_n | \beta)$$

$$= \operatorname{argmax}_{\beta} L(\beta) \quad \text{Due to monotonicity of log transform}$$

$$= \operatorname{argmax}_{\beta} \log L(\beta)$$

$$= \operatorname{argmax}_{\beta} \log \left\{ \left(\frac{1}{2}\right)^n \prod_{i=1}^n e^{-|y_i - \beta^T x_i|} \right\}$$

$$= \operatorname{argmax}_{\beta} \left\{ \log \left[\underbrace{\left(\frac{1}{2}\right)^n}_{\text{doesn't depend on } \beta} + \log \left[\prod_{i=1}^n e^{-|y_i - \beta^T x_i|} \right] \right] \right\}$$

$$= \operatorname{argmax}_{\beta} \left\{ \log \prod_{i=1}^n e^{-|y_i - \beta^T x_i|} \right\}$$

$$= \operatorname{argmax}_{\beta} \left\{ \sum_{i=1}^n \log e^{-|y_i - \beta^T x_i|} \right\}$$

Since $\log e^{-x} = -x$

$$= \operatorname{argmax}_{\beta} \left\{ \sum_{i=1}^n -|y_i - \beta^T x_i| \right\}$$

$$= \operatorname{argmax}_{\beta} \left\{ - \sum_{i=1}^n |y_i - \beta^T x_i| \right\} \quad \text{Since } \operatorname{argmax}_{\beta} \{-f(\beta)\} =$$

$$\operatorname{argmin}_{\beta} \{f(\beta)\}$$

$$= \operatorname{argmin}_{\beta} \{|y_i - \beta^T x_i|\} \quad \text{Hence Proved}$$

Problem 4

a) We have here the union of 3 disjoint intervals

We must show that $\{1, 2, 3, 4, 5, 6\}$ is shattered

Suppose the labellings are $(c_1, c_2, c_3, c_4, c_5, c_6)$
where c_i is the labelling of i

b) We must show that all 2^6 labellings are possible

Let us associate $[a_1, b_1] \rightarrow c_1, c_2$ for 1, 2

$[a_2, b_2] \rightarrow c_3, c_4$ for 3, 4

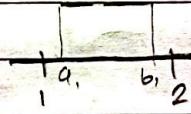
$[a_3, b_3] \rightarrow c_5, c_6$ for 5, 6

WLOG, if one interval can produce all possible 2^2 labellings
for 2 points, then the union of 3 disjoint intervals can
generate $(2^2)^3 = 2^6$ labellings

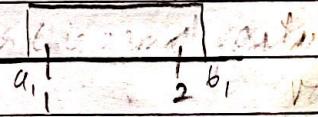
c) WLOG, let's take $[a_1, b_1]$ & try to shatter $\{1, 2\}$

Let ϵ be a small +ve value that is $0 < \epsilon \ll 1$

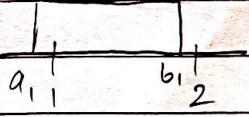
Labelling: 00 \rightarrow if $a_1 = -1 + \epsilon$
 $b_1 = 2 - \epsilon$



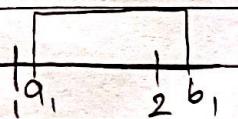
d) Labelling: 11 \rightarrow if $a_1 = 1 - \epsilon$
 $b_1 = 2 + \epsilon$



e) Labelling: 10 \rightarrow if $a_1 = 1 - \epsilon$
 $b_1 = 2 - \epsilon$



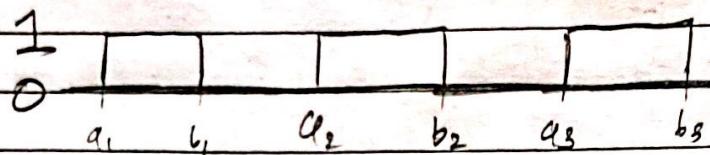
f) Labelling: 01 \rightarrow if $a_1 = 1 + \epsilon$
 $b_1 = 2 + \epsilon$



∴ Since 2^2 labellings are possible for 1 interval $(2^2)^3$ is possible for 3
 $\{1, 2, 3, 4, 5, 6\}$ is shattered

b) Suppose we have some sample of size 7, say $\{1, 2, 3, 4, 5, 6, 7\}$

We know that the graph of the hypothesis is as follows:



8 There are 3 rising edges & 3 dipping edges.

For a sample of 7 points, consider the labelling

1 0 1 0 1 0 . 1

This particular labelling has 4 separate rises, & we would need 4 disjoint intervals to generate this.

However, this function only gives us 3 jumps but we need 4 to generate one of the labellings, specifically 1010101

\therefore A set of size 7 may not be shattered

This is because the function class of union of 3 disjoint intervals is not sufficiently complex.