

Connor Relalmu

I please my horses
that have abdicated
the Stevens theorem
System confirmation

1.

Class $\rightarrow w_1, w_2 \dots w_c$

a.) Bayes Theorem: $P(w_i|x) = \frac{P(x|w_i)P(w_i)}{P(x)} \quad i=1, \dots, c$

$$\text{where } P(x) = \sum_{i=1}^c P(x|w_i)P(w_i)$$

Rule: choose w_i over w_j (where $i, j \in \{1, \dots, c\}$ & $i \neq j$) if
the posterior probability of w_i is greater than w_j .

i.e. if $P(w_i|x) > P(w_j|x)$, choose w_i .

Using only the class conditional density & prior would yield:

$$\frac{P(x|w_i)P(w_i)}{P(x)} \rightarrow \frac{P(x|w_i)P(w_i)}{P(x)}$$

Since $P(x)$ is the same for both it can be ignored.

\therefore decide w_i if:

likelihood ratio
 \downarrow

$$P(x|w_i)P(w_i) > P(x|w_j)P(w_j)$$

or

$$\frac{P(x|w_i)}{P(x|w_j)} > \frac{P(w_i)}{P(w_j)}$$

where $i, j \in \{1, \dots, c\}$
& $i \neq j$

$$b.) \quad w_1 \quad p(x|w_1) = \frac{1}{\sqrt{2\pi}} e^{-\left(\frac{x^2}{2}\right)}$$

$$p(x|w_2) = 1, x \in [-\frac{1}{2}, \frac{1}{2}]$$

w/ loss taken into respect we decide based
on minimizing the risk:

$$R(a_i|x) = \sum_{j=1}^J \lambda_{ij} p(w_j|x)$$

where λ_{ij} is the loss for taking action a_i (choose w_i)

when the true state of nature is w_j

w/ 2 classes

$$R(a_1|x) = \lambda_{11} p(w_1|x) + \lambda_{12} p(w_2|x)$$

$$R(a_2|x) = \lambda_{21} p(w_1|x) + \lambda_{22} p(w_2|x)$$

$$\text{w/ zero 1 loss} \rightarrow \lambda = [0 \ 0]$$

$$R(a_1|x) = p(w_2|x)$$

$$R(a_2|x) = p(w_1|x) \rightarrow \text{Decision rule} \rightarrow \begin{array}{l} \text{choose } w_1 \text{ if } \\ p(w_2|x) < p(w_1|x) \end{array}$$

minimize risk

This is the same result as part a.) i.e. we can use:

$$\text{Choose } w_1 \rightarrow \frac{p(x|w_1)}{p(x|w_2)} > \frac{p(w_2)}{p(w_1)} \rightarrow p(x|w_1) > p(x|w_2)$$

Since $p(x|w_1)$ is > 0 for all $x \rightarrow p(x|w_1) > p(x|w_2)$ when
 x is not in the interval $[-\frac{1}{2}, \frac{1}{2}]$



Otherwise, choose w_1 if

$$P(x|w_1) > 1$$

$$\frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} > 1 \rightarrow -\frac{x^2}{2} > \ln(\sqrt{2\pi}) \rightarrow -x^2 > 2\ln(\sqrt{2\pi})$$

$$\rightarrow x > (-2\ln(\sqrt{2\pi}))^{\frac{1}{2}} \rightarrow P(x|w_1) \text{ is never equal to 1.}$$

∴ Overall decision rule:

$$\text{Rule : } \begin{cases} \text{Choose } w_2, x \in [-\frac{1}{2}, \frac{1}{2}] \\ \text{Choose } w_1, \text{ Otherwise} \end{cases}$$

2. a.)

$$E_p(w) = \sum_{i=1}^N -\{y_i \log f(x_i) + (-y_i) \log (1-f(x_i))\} + \frac{\lambda}{2} \|w\|^2$$

$$\Delta_w E_p(w) = \sum_{i=1}^N -\frac{d}{dw} (y_i \log \delta(w^\top x_i) + (-y_i) \log (1-\delta(w^\top x_i))) + \frac{d}{dw} \frac{\lambda}{2} \|w\|^2$$

$$= - \left[\frac{y_i}{\delta(w^\top x_i)} \left(\frac{d\delta(w^\top x_i)}{dw} \right) + \frac{-y_i}{1-\delta(w^\top x_i)} \left(\frac{-d\delta(w^\top x_i)}{dw} \right) \right]$$

$$= - \left[\frac{y_i}{\delta(w^\top x_i)} - \frac{-y_i}{1-\delta(w^\top x_i)} \right] \frac{d(\delta(w^\top x_i))}{dw}$$

$$= - \left[\frac{y_i}{\delta(w^\top x_i)} - \frac{-y_i}{1-\delta(w^\top x_i)} \right] \delta(w^\top x_i) (1-\delta(w^\top x_i)) x_i$$

$$= - \left[y_i (1-\delta(w^\top x_i)) - (-y_i) (\delta(w^\top x_i)) \right] x_i$$

$$= - \left[y_i - \delta(w^\top x_i) x_i \right]$$

$$= (\delta(w^\top x_i) - y_i) x_i$$

$$\frac{d}{dw} \frac{\lambda}{2} \|w\|^2 = \lambda w$$

i)

$$\Delta_w E_\rho(w) = \sum_{i=1}^N (\delta(w^\top x_i) - y_i)x_i + \lambda w$$

b.)

$$P(Y|w) = \prod_{i=1}^N \delta(w^\top x_i)^{y_i} (1 - \delta(w^\top x_i))^{(1-y_i)}$$
$$P(w) = ?$$

$P(Y|w)$ is equivalent to the above calculation w/o the negative.
So if $\Delta_w \ln(P(w)) = \lambda w$ then you would get the same results as in part A
If you maximize $P(Y|w) P(w)$,



(I was spending time on part 3 because I thought I needed a converged result, I ran out of time to look back at this)

3.)

| Instance label | (1) | (2) | (3) | (4) | (5) | (6) | (7) | (8) |
|-------------------|---------|-------|-------|-------|-------|------------|-------|-------|
| Data (x_1, x_2) | (10,10) | (0,0) | (8,4) | (3,3) | (4,8) | (0.5, 0.5) | (4,3) | (2,5) |

Positive $\rightarrow t = 1$, negative $t = -1$

Start at $x_1 + x_2 = 0$ i.e. $w = [0 \ 1 \ 1]^\top$ $x = [x_1 \ x_2]^\top$
 $t_i w^\top x_i > 0 \rightarrow \text{classification correct}$

$t=0 \quad P(1) \quad x_1 = 10, x_2 = 10, t_1 = 1$

$$t_1 w^\top x_1 = 1 [0 \ 1 \ 1] \begin{bmatrix} 1 \\ 10 \\ 10 \end{bmatrix} = (0+10) = 20$$

$$t_1 w^\top x_1 > 0 \rightarrow \text{keep } w \rightarrow w^1 = w^0$$

$t=1 \quad P(2) \quad x_1 = 0, x_2 = 0, t_2 = -1$

$$t_2 w^\top x_2 = -1 [0 \ 1 \ 1] \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = 0$$

$t_2 w^\top x_2 \leq 0 \rightarrow \text{update } w$

$$w^2 = w^1 + t_2 x_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix}$$

$t=2 \quad P(3) \quad x_1 = 8, x_2 = 4, t_3 = 1$

$$t_3 w^\top x_3 = 1 [-1 \ 1 \ 1] \begin{bmatrix} -1 \\ 8 \\ 4 \end{bmatrix} = -1 + 8 + 4 = 11$$

$t_3 w^\top x_3 > 0 \rightarrow \text{keep } w \rightarrow w^3 = w^2$

$$T=3 \quad P(4) \quad x_1 = 3, x_2 = 3, t_4 = -1$$

$$t_4 w^T x_4 = -1 \begin{bmatrix} -1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 3 \end{bmatrix} = -(-1+3+3) = -5$$

$t_4 w^T x_4 \leq 0 \rightarrow \text{update } w$

$$w^4 = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 3 \\ 3 \end{bmatrix} = \begin{bmatrix} -2 \\ -2 \\ -2 \end{bmatrix}$$

$$T=4 \quad P(5) \quad x_1 = 4, x_2 = 8, t_5 = 1$$

$$t_5 w^T x_5 = 1 \begin{bmatrix} -2 & -2 & -2 \end{bmatrix} \begin{bmatrix} 1 \\ 4 \\ 8 \end{bmatrix} \leftarrow 0 \rightarrow \text{update } w$$

$$w^5 = \begin{bmatrix} -2 \\ -2 \\ -2 \end{bmatrix} + \begin{bmatrix} 1 \\ 4 \\ 8 \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \\ 6 \end{bmatrix}$$

$$T=5 \quad P(6) \quad x_1 = 4, x_2 = 4, t_6 = -1$$

$$t_6 w^T x_6 = -1 \begin{bmatrix} -1 & 2 & 6 \end{bmatrix} \begin{bmatrix} 1 \\ 4 \\ 4 \end{bmatrix} = -1(-1+8+24) \leftarrow 0 \rightarrow \text{update } w$$

$$w^6 = \begin{bmatrix} -1 \\ 2 \\ 6 \end{bmatrix} - \begin{bmatrix} 1 \\ 4 \\ 4 \end{bmatrix} = \begin{bmatrix} -2 \\ 15 \\ 5 \end{bmatrix}$$

$$T=6 \quad P(7) \quad x_1 = 4, x_2 = 3, t_7 = 1$$

$$t_7 w^T x_7 = 1 \begin{bmatrix} -2 & 1 & 5 & 1 & 5 \end{bmatrix} \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} = -2+6+16 \geq 0 \rightarrow w^8 = w^7$$

$$\gamma = \rho(g) \quad x_1 = 2 \quad x_2 = 5 \quad t_g = 1$$

$$t_g w^T x_g - 1 \begin{bmatrix} -2, 5, 5, 5 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 5 \end{bmatrix} > 0 \Rightarrow w^9 = w^8$$

$$\text{Final result} \rightarrow w = \begin{bmatrix} -2, 5, 5, 5 \end{bmatrix}^T$$

U_i

a.) $p(x) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n k\left(\frac{x-x_i}{h}\right)$

$$k(u) = \begin{cases} 1, & |u| \leq 1/2 \\ 0, & \text{o.w.} \end{cases}$$

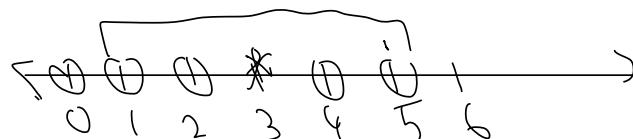
$$p(3) = \frac{1}{5} \sum_{i=1}^5 \frac{1}{5} k\left(\frac{3-x_i}{5}\right) \quad x_i \in \{0, 1, 2, 4, 5\}$$

$$\frac{1}{25} \left(k\left(\frac{3}{5}\right) + k\left(\frac{2}{5}\right) + k\left(\frac{1}{5}\right) + k\left(-\frac{1}{5}\right) + k\left(-\frac{2}{5}\right) \right)$$

\uparrow \uparrow \uparrow \uparrow \uparrow
 $x_1 = 0$ $x_2 = 1$ $x_3 = 1$ $x_4 = -1$ $x_5 = -2$

$$p(3) = \frac{1}{25} (4) = \frac{4}{25}$$

4 total around 3 w/ n=5



b) used Euclidean distance we can visually
see which points are closest to \mathbf{x}

i.) w/ $k=1 \rightarrow 1\text{NN} \rightarrow \mathbf{x}$ is from red class

$$\hookrightarrow \frac{k_{\text{Red}}}{k} = \frac{1}{14}, \frac{k_{\text{Blue}}}{k} = \frac{0}{14} \quad \frac{k_{\text{Red}}}{k} > \frac{k_{\text{Blue}}}{k}$$

ii.) w/ $k=7 \rightarrow 7\text{NN} \rightarrow \mathbf{x}$ is from blue class

$$\hookrightarrow \frac{k_{\text{Red}}}{k} = \frac{3}{14}, \frac{k_{\text{Blue}}}{k} = \frac{4}{14}, \quad \frac{k_{\text{Blue}}}{k} > \frac{k_{\text{Red}}}{k}$$

iii.) I don't believe we should use $k=9$ since
the other training samples are relatively far
away from the test sample when compared
to a 1NN or 7NN classifier. If $k=9$,
the degree of smoothing is likely to be too large
and the classifier could become inaccurate.

$$5.) f(x|\theta) = 2\theta x e^{-\theta x^2} \quad x > 0 \quad x \rightarrow (x_1, x_2, \dots, x_n)$$

\nwarrow Likelihood

$$f(x|\theta) = \prod_{i=1}^n 2\theta x_i e^{-\theta x_i^2}$$

$$\text{Log-Likelihood} \rightarrow \ln(f(x|\theta)) = \sum_{i=1}^n \ln(2\theta x_i e^{-\theta x_i^2})$$

$$LL(\theta) = \sum_{i=1}^n (\ln(2\theta x_i) - \theta x_i^2)$$

$$\Delta_\theta LL(\theta) = \sum_{i=1}^n \left(\frac{2x_i}{2\theta x_i} - x_i^2 \right) = \sum_{i=1}^n \left(\frac{1}{\theta} - x_i^2 \right) = \frac{n}{\theta} - \sum_{i=1}^n x_i^2$$

$$\text{Set } \frac{\partial}{\partial \theta} \sum_{i=1}^n x_i^2 = 0 \rightarrow \frac{1}{\theta} = \frac{1}{n} \sum_{i=1}^n x_i^2$$

$$\text{MLE} \rightarrow \theta = \frac{n}{\sum_{i=1}^n x_i^2} \rightarrow \frac{1}{\theta} = \frac{1}{n} \sum_{i=1}^n x_i^2$$

$$\text{b) } \text{Unif}(\theta_1, \theta_2) \quad \left\{ \begin{array}{l} \frac{1}{\theta_2 - \theta_1}, x_i \text{ between } \theta_1, \theta_2 \text{ for } x_i \in \{x_1, \dots, x_n\} \\ 0, \text{ o.w.} \end{array} \right.$$



$$X_i \sim \text{Unif}(\theta_1, \theta_2) = \prod_{i=1}^n \frac{1}{\theta_2 - \theta_1} X_i \text{ iid } \theta_1, \theta_2 \text{ O.w.}$$

$$\text{Likelihood} \rightarrow = \left(\frac{1}{\theta_2 - \theta_1} \right)^n$$

$$LL(\theta) \rightarrow \ln \left(\frac{1}{\theta_2 - \theta_1} \right)^n = -n \ln (\theta_2 - \theta_1)$$

$$\Delta_{\theta_1} LL(\theta) = -n \left(\frac{1}{\theta_2 - \theta_1} \right) \frac{d}{d\theta_1} (\theta_2 - \theta_1) = \frac{n}{\theta_2 - \theta_1}$$

$$\Delta_{\theta_2} LL(\theta) = -n \left(\frac{1}{\theta_2 - \theta_1} \right) \frac{d}{d\theta_2} (\theta_2 - \theta_1) = \frac{-n}{\theta_2 - \theta_1}$$

for θ_1 , $\Delta_{\theta_1} LL(\theta) = \frac{n}{\theta_2 - \theta_1} \rightarrow$ since $X_i \sim \text{Unif}[\theta_1, \theta_2]$
 we maximize θ_1 by
 choosing the smallest X value
 $\theta_1 = \min(X) \quad X \in \{x_1, x_2, \dots, x_n\}$

for θ_2 , $\Delta_{\theta_2} LL(\theta) = \frac{-n}{\theta_2 - \theta_1} /$ to maximize θ_2 , we
 need to choose the largest
 X value, $\theta_2 = \max(X) \quad X \in \{x_1, x_2, \dots, x_n\}$

6.)

a.) The discriminative approach makes decisions by directly modeling the posterior probability of the classes, i.e. $P(C_k|x)$. The generative

and uses Bayes' rule to find posterior.

approach focuses on modeling the class conditional densities $P(x|C_k)$ along with the priors $P(C_k)$. The discriminative function approach utilizes the creation of a function, $g(x)$ which maps a sample to a particular class/label

b.) When performing MLE, the assumption is that the parameters are a fixed variable,

while Bayesian estimation assumes the parameter comes from an underlying distribution. This makes MLE computationally efficient while Bayes can possibly allow for more accurate results if prior information is reliable.

c.) PCA & FLD are fairly similar techniques but with different goals. PCA is an unsupervised technique that is aiming to reduce dimensionality without losing information. FLD is supervised and requires the class labels. This allows FLD to be used to find the best separation across

classes.

- d.) Batch gradient descent is performed using all the training data at once, while mini-batch iterates through small portions of the data to perform updates. Stochastic gradient descent iterates through samples one at a time to make updates.

Batched GD is generally unbiased, accurate, and can converge to a global optimum. However, it is slow and unreliable w/ large datasets & can easily fall into traps w/ local modes.