

ECE523: Engineering Applications of Machine
Learning and Data Analytics

Due 01/26/2018 @ 11:59PM (D2L)

I acknowledge that this homework is solely my effort. I have done this work by myself. I have not consulted with others about in any way. I have not received outside aid (outside of my own brain). I understand that violation of these rules contradicts the class policy on academic integrity.

Name:

David Akre

Signature:

David Akre

Date:

1/24/18

Instructions: There are seven problems. Partial credit is given for answers that are partially correct. No credit is given for answers that are wrong or illegible. All work must be supported and code must be submitted for credit.

Theory: _____

Practice: _____

Total: _____

① Maximum Posterior vs. Probability of chance

- (a) Show/Explain that $p(w_{\max}|x) \geq \frac{1}{c}$ when we are using the Bayes Decision Rule. Derive an expression for $p(\text{error})$ - let w_{\max} be the true state of nature for which $p(w_{\max}|x) \geq p(w_i|x)$ for $i=0,1,\dots,c$. Show that $p(\text{error}) \leq (c-1)/c$ when we use the Bayes decision rule to make a decision.

Step 1: Derive an expression for $p(\text{error})$

$$\begin{aligned}
 p(\text{error}) &= \int_1^c p(\text{error}, x) dx \quad \leftarrow \therefore p(\text{error}, x) = p(\text{error}|x) p(x) \\
 &= \int_1^c p(\text{error}|x) p(x) dx \quad \leftarrow \text{Choose class with highest posterior (i.e. } 1 - p(w_{\max}|x)) \\
 &= \int_1^c (1 - p(w_{\max}|x)) dx \\
 &= \int_1^c dx - \int_1^c p(w_{\max}|x) dx \\
 &= x \Big|_1^c - p(w_{\max}|x) x \Big|_1^c \\
 &= x(1 - p(w_{\max}|x)) \Big|_1^c \\
 &= c(1 - p(w_{\max}|x)) - (1 - p(w_{\max}|x)) \\
 &= c - c p(w_{\max}|x) - 1 + p(w_{\max}|x) \quad \leftarrow \therefore \text{Divide by } c \text{ and cancel } p(w_{\max}|x) \\
 &= \frac{c-1}{c}
 \end{aligned}$$

$$\boxed{\therefore p(\text{error}) \leq \frac{c-1}{c}}$$

Step 2: Show that $p(w_{\max}|x) \geq \frac{1}{c}$ for Bayes Decision Rule

$$p(w_{\max}|x) = \frac{p(x|w_{\max}) p(w_{\max})}{p(x)} \quad \leftarrow p(x|w_{\max}) \text{ is the likelihood probability, } p(w_{\max}) \text{ is the prior probability, and } p(x) \text{ is the conditional or evidence}$$

$\therefore w_{\max}$ resembles the class with the highest posterior probability which in turn means that it resembles the minimum probability of error. So for $i=0,1,\dots,c$ $p(w_{\max}|x) \geq p(w_i|x)$ means that $p(\text{error}) \leq \frac{c-1}{c}$.

$$1 - p(\text{error}|x) \geq p(w_i|x)$$

$$1 \geq p(\text{error}|x) + p(w_i|x) \quad \leftarrow \text{NOTE: The error summed with } p(w_i|x) = 1$$

$$1 \geq p(\text{error}|x) + \sum_{i=1}^c p(w_i|x) \quad \leftarrow \lim_{i \rightarrow 0} \sum_{i=1}^c p(w_i|x) = 1 - p(w_{\max}|x) \text{ or } \frac{1}{c}$$

NOTE: The posterior probability must be greater than or equal to the inverse of the number of classes, otherwise a different class is the real w_{\max} (e.g. $\sum_{i=0}^c p(w_i|x) \leq p(w_{\max}|x)$)

$$1 \geq 1 - p(w_{\max}|x) + \frac{1}{c}$$

$$-\frac{1}{c} \geq -p(w_{\max}|x)$$

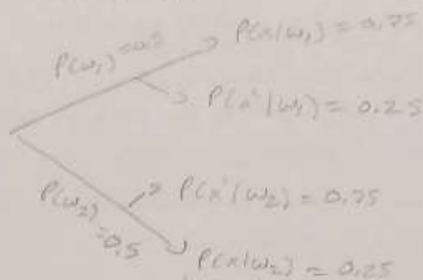
$$\boxed{\frac{1}{c} \leq p(w_{\max}|x)}$$

Step 3: Show proof in an example

(i) Two classes: w_1 & w_2

- $P(w_1) = 0.5$
- $P(w_2) = 0.5$
- $P(x|w_1) = 0.75$
- $P(x|w_2) = 0.25$

(ii) Decision Tree



(iii) Total probability (evidence) $P(x) = P(w_1)P(x|w_1) + P(w_2)P(x|w_2)$
 $= (0.5)(0.75) + (0.5)(0.25)$
 $= 0.5 \rightarrow P(x) = 0.5$

(iv) Distribution for all events represented by the following table:

$P(x w_1)P(w_1) = 0.375$
$P(x' w_1)P(w_1) = 0.125$
$P(x w_2)P(w_2) = 0.125$
$P(x' w_2)P(w_2) = 0.375$

(v) Using Bayes theorem find $P(w_1|x)$ & $P(w_2|x)$:

$$P(w_1|x) = \frac{P(x|w_1)P(w_1)}{P(x)} = 0.75$$

$$P(w_2|x) = \frac{P(x|w_2)P(w_2)}{P(x)} = 0.25$$

(vi) Find w^* using Bayes Decision Rule:

- $P(w_1|x) \geq P(w_2|x)$; so $P(w_1|x)$ is selected since it has a higher posterior probability

\therefore Check $P(w_1|x) \geq \frac{1}{c}$ where $c=2$

$$0.75 \geq 0.5 \leftarrow \text{True (first proof holds)}$$

(vii) Verify $p(\text{error}) \leq \frac{c}{c-1}$

$$\begin{aligned} - p(\text{error}|x) &= 1 - p(w_{\max}|x) ; w_{\max} = w_1 & - p(\text{error}|x) &= 1 - p(w_2|x) \\ &= 1 - 0.75 & &= 1 - 0.25 \\ &= 0.25 & &= 0.75 \end{aligned}$$

$\therefore 0.25 < 0.75$ (w_1 leads to a smaller probability of error)

$$- 0.25 \leq \frac{2-1}{2}$$

$$0.25 \leq 0.5 \leftarrow \text{True (second proof holds)}$$

② Bayes Decision Rule Classifier

(a) Let the elements of a vector $\mathbf{x} = [x_1, \dots, x_d]^T$ be binary values, but $P(w_j)$ be the prior probability of the class w_j ($j \in C = \{1, 2\}$) and let $p_{ij} = P(x_i = 1 | w_j)$ with all elements in \mathbf{x} being independent. If $P(w_1) = P(w_2) = \frac{1}{2}$, and $p_{i1} = p > \frac{1}{2}$ and $p_{i2} = 1 - p$ show that the minimum error decision rule is: choose w_1 if $\sum_{i=1}^d x_i > \frac{d}{2}$. HINT: choose w_1 if $P(w_1)P(\mathbf{x}|w_1) > P(w_2)P(\mathbf{x}|w_2)$.

Step 0: Describe given situation as a problem

- Given: $\mathbf{x} = [x_1, \dots, x_d]^T$ — $\mathbf{x} \in \{0, 1\}^d$
- $\mathcal{C} = \{w_1, w_2\}$ (Two classes) — $1 \in \mathcal{C} \subseteq \mathcal{D}$
- $P(w_1) = P(w_2) = \frac{1}{2}$ — $x_i \in \{0, 1\}$ ← Binary (Bernoulli)
- $p_{ij} = P(x_i = 1 | w_j)$

Prove: choose w_1 if $\sum_{i=1}^d x_i > \frac{d}{2}$

∴ Bayes Decision Rule: choose class with highest posterior probability which leads to the min error probability

Step 1: Start with Bayes Decision Rule to simplify down to a condition

— Select the maximum posterior to find the minimum probability of error

$$\max \{ P(w_1)P(\mathbf{x}|\mathbf{w}_1), P(w_2)P(\mathbf{x}|\mathbf{w}_2) \}$$

— We want to prove w_1 is max posterior if $\sum_{i=1}^d x_i > \frac{d}{2}$

$$P(w_1)P(\mathbf{x}|\mathbf{w}_1) > P(w_2)P(\mathbf{x}|\mathbf{w}_2)$$

— We're also told that $P(w_1) = P(w_2) = \frac{1}{2}$

$$P(\mathbf{x}|\mathbf{w}_1) > P(\mathbf{x}|\mathbf{w}_2)$$

$$P(\mathbf{x}|\mathbf{w}_1)$$

$$P(\mathbf{x}|\mathbf{w}_1) > P(\mathbf{x}|\mathbf{w}_2) \leftarrow \text{Condition we want}$$

Step 2: Utilizing independence rule write expression out as product (naive Bayes)
 \hookrightarrow joint model

$$\prod_{i=1}^d P(x_i | w_1) > \prod_{i=1}^d P(x_i | w_2)$$

Step 3: From the joint model, rewrite each distribution over the class variable w

$$(i) \left(\prod_{i=1}^d P(x_i | w_1) \right) \frac{1}{2} P(w_1) = P(w_1 | \mathbf{x})$$

$$Z = P(\mathbf{x}) = \sum_{i=1}^d P(w_1)P(\mathbf{x}|\mathbf{w}_1)$$

$$= \sum_{i=1}^d P(x_i | w_1) \leftarrow \text{NOTE: same simplification applied to } w_2$$

$$(ii) \sum_{i=1}^d P(x_i | w_1) > \sum_{i=1}^d P(x_i | w_2)$$

Step 4: Utilizing the last hint to prove " w_1 if $\sum_{i=1}^d x_i > \frac{d}{2}$ "

$$\text{Given } p_{i1} = p > \frac{1}{2}$$

$$p_{i2} = 1 - p$$

→ Find the posterior for w_1 & w_2 (analogous for w_2)

$$P_i(x) = \sum_{j=1}^2 P(x_i | w_j) P(w_j)$$

$$\frac{P(x | w_1) P(w_1)}{P(w_1 | x)} = \sum_{i=1}^d \frac{P(x_i | w_1) P(w_1)}{P(w_1 | x)}$$

$$(i) P(w_1 | x) = \frac{P(x | w_1) P(w_1)}{\sum_{i=1}^d P(x_i | w_1) P(w_1)} > \frac{1}{2}$$

$$(ii) P(w_2 | x) = \frac{P(x | w_2) P(w_2)}{\sum_{i=1}^d P(x_i | w_2) P(w_2)} = 1 - P$$

∴ If $P(w_1 | x) > \frac{1}{2}$ then $\sum_{i=1}^d x_i > \frac{d}{2}$ because if the running sum of 1's in the Bernoulli distribution is less than or equal to $d/2$ then $P(w_1 | x_i = 1)$ cannot be the maximum posterior which leads to the minimum probability of error. If it is less than or equal to $d/2$ then there are more 0's in the Bernoulli set which results in w_2 being the posterior max.

3) The Ditzler Household Growing up

My parents have two kids now grown into adults. Obviously, then I am a boy. I was born on a Weds. What is the probability that I have a brother? You can assume $P(\text{boy}) = P(\text{girl}) = \frac{1}{2}$

Step 1: Calculate prior probabilities

- Given: $P(\text{boy}) = P(\text{girl}) = \frac{1}{2}$

- Priors: boy & boy - $P(bb) = (\frac{1}{2})(\frac{1}{2}) = \frac{1}{4}$
 boy & girl - $P(bg) = (\frac{1}{2})(\frac{1}{2}) = \frac{1}{4}$
 girl & boy - $P(gb) = (\frac{1}{2})(\frac{1}{2}) = \frac{1}{4}$
 girl & girl - $P(gg) = (\frac{1}{2})(\frac{1}{2}) = \frac{1}{4}$ } $b_g + g_b = \frac{1}{2}$
 Not possible

Step 2: Calculate probability there are two boys given one boy already in the family:

$$P(bb|b) = \frac{P(b|bb)P(bb)}{P(b)} = \frac{(1)(0.25)}{(0.75)} = \frac{1}{3}$$

- $P(b|bb)$ = Probability that both are boys given one is already a boy (1)

- $P(b)$ = Probability of a boy ($\frac{1}{2} + \frac{1}{4}$)

$$P(bb|b) = \frac{1}{3} \text{ or } 0.33 \text{ or } 33\%$$

4) Linear Classifier with a Margin

9) Show that regardless of the dimensionality of the feature vectors, a data set that has just two data points, one from each class, is sufficient to determine the location of the maximum hyperplane.
 HWT #1: Consider a data set of two data points, $x_1 \in C_1$ ($y_1 = +1$) and $x_2 \in C_2$ ($y_2 = -1$) and set up the minimization problem (for computing the hyperplane) with appropriate constraints on $w^T x_1 + b$ and $w^T x_2 + b$ and solve it. HWT #2: This can be framed as a constrained optimization problem $\arg \min_{w \in \mathbb{R}^d} \|w\|_2^2$ (subject to some constraint). What is w & b ? What are the constraints? How

Did we solve the constrained optimization problem in Fisher's Linear Discriminate?

Step 0: Lay out given properties

- Data set $\rightarrow x_1 \in C_1$ ($y_1 = +1$) - Constraints $\rightarrow w^T x_1 + b = 1$ - optimization $\rightarrow \min \|w\|_2^2$ (minimize square norm)
 $x_2 \in C_2$ ($y_2 = -1$) $\rightarrow w^T x_2 + b = -1$

\therefore We want to maximize the margin between the two data points

Step 1: Apply Lagrange multipliers to combined constraints

$$y_1(w^T x_1 + b) = 1$$

$$y_2(w^T x_2 + b) = 1 \rightarrow (\text{same as } -y_2(w^T x_2 + b) = -1)$$

- Recasting constraints as functions of $g_1(x_1, y_1)$ & $g_2(x_2, y_2)$

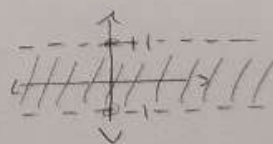
$$g_1(x_1, y_1) = y_1(w^T x_1 + b) - 1$$

$$g_2(x_2, y_2) = y_2(w^T x_2 + b) + 1$$

- write out Lagrange multipliers

$$d_1(x_1, y_1, \lambda) = f(x_1, y_1) + \lambda \cdot g_1(x_1, y_1)$$

$$d_2(x_2, y_2, \lambda) = f(x_2, y_2) + \lambda \cdot g_2(x_2, y_2)$$



$$L_1(x_1, y_1, \lambda) = \arg \min_{x \in \mathbb{R}^2} \|w\|_2^2 + (y_1 (w^T x_1 + b) - 1) \lambda$$

$$L_2(x_2, y_2, \lambda) = \arg \min_{x \in \mathbb{R}^2} \|w\|_2^2 + (y_2 (w^T x_2 + b) + 1) \lambda$$

- Take the derivative of each Lagrangian function wrt λ and set it to 0

$$\begin{aligned} \text{(i)} \quad \frac{dL_1}{d\lambda} &= 0 + y_1 w^T x_1 + y_1 b - 1 = 0 \Rightarrow \underline{y_1 (w^T x_1 + b) - 1 = 0} \\ \text{(ii)} \quad \frac{dL_2}{d\lambda} &= 0 + y_2 w^T x_2 + y_2 b + 1 = 0 \Rightarrow \underline{y_2 (w^T x_2 + b) + 1 = 0} \end{aligned} \quad \left. \vphantom{\begin{aligned} \text{(i)} \quad \frac{dL_1}{d\lambda} = 0 + y_1 w^T x_1 + y_1 b - 1 = 0 \Rightarrow \underline{y_1 (w^T x_1 + b) - 1 = 0} \\ \text{(ii)} \quad \frac{dL_2}{d\lambda} = 0 + y_2 w^T x_2 + y_2 b + 1 = 0 \Rightarrow \underline{y_2 (w^T x_2 + b) + 1 = 0} \right\} \text{original constraints}$$

- Take derivative of each Lagrangian function wrt x or wrt y and set it to 0

$$\text{(iii)} \quad \frac{dL_1}{dx_1} = 0 + y_1 w^T \lambda = 0 \Rightarrow \underline{y_1 w^T \lambda = 0} ; \underline{y_1 = 0}$$

$$\text{(iv)} \quad \frac{dL_2}{dx_2} = 0 + y_2 w^T \lambda = 0 \Rightarrow \underline{y_2 w^T \lambda = 0} ; \underline{y_2 = 0}$$

$$\text{(v)} \quad \frac{dL_1}{dy_1} = 0 + w^T x_1 \lambda + b \lambda = 0 \Rightarrow \underline{\lambda w^T x_1 + b = 0} ; \underline{x_1 = \frac{-b\lambda}{w^T \lambda} = \frac{-b}{w^T}}$$

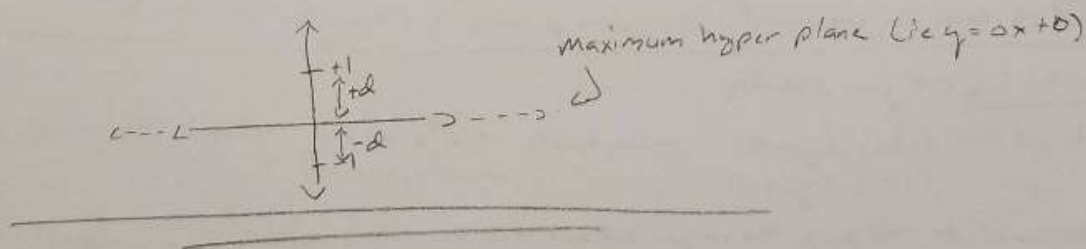
$$\text{(vi)} \quad \frac{dL_2}{dy_2} = 0 + w^T x_2 \lambda + b \lambda = 0 \Rightarrow \underline{\lambda w^T x_2 + b = 0} ; \underline{x_2 = \frac{-b\lambda}{w^T \lambda} = \frac{-b}{w^T}}$$

- Substitute x_1, x_2, b_1 , and y_2 back into (i), (ii)

$$\begin{aligned} &0 (w^T (\frac{-b}{w^T}) + b) - 1 = 0 \\ &0 - 1 = 0 \quad \therefore \text{Nothing} = -1 \text{ which is nothing (i.e. 0)} \end{aligned}$$

$$\begin{aligned} &0 (w^T (\frac{-b}{w^T}) + b) + 1 = 0 \\ &0 + 1 = 0 \quad \therefore \text{Nothing} = 1 \text{ which is nothing (i.e. 0)} \end{aligned}$$

\therefore Thus the Lagrange multipliers end up being 0 (i.e. $\lambda = 0$). This means the maximum hyperplane lies on the origin at 0,0 which produces the largest distances from the two data points +1 & -1



⑤ Decision making with Bayes

(a) The Bayes Decision Rule describes the approach we take choosing a class c for a data point x . This can be achieved modeling $P(w|x)$ or $P(x|w)P(w)/P(x)$. Compare and contrast these two approaches to modeling and discuss the advantages and disadvantages. For the latter model, why might learning $P(x)$ be useful?

Part 1: $P(w|x)$

- $P(w|x)$ describes the posterior probability that some event will occur given a second event also occurs. Since we're talking about Bayes Decision Rule, $P(w|x)$ represents the maximum posterior probability (aka $\max P(w|x)$).
- Advantages: of this model is that it describes the Bayes decision rule in a more simplistic way (e.g. the probability w occurs given x also occurs).
- Disadvantages: This model does not describe in detail the likelihood, prior, or evidence probabilities (e.g. $P(x|w)$, $P(w)$, $P(x)$).
- Compare/contrast: $P(w|x)$ can be expanded out to describe $P(x|w)P(w)/P(x)$, so overall they represent the same results. The contrasting point is how they're expressed (e.g. $P(x|w)P(w)/P(x)$ shows likelihood \times prior / evidence, whereas $P(w|x)$ just shows the posterior probability).

Part 2: $P(x|w)P(w)/P(x)$

- $P(x|w)P(w)/P(x)$ describes the combination of the sum & product rules to describe Bayes. It explicitly shows the likelihood, prior, and evidence probabilities.
- Advantage: $P(x|w)P(w)/P(x)$ is a more verbose way to express the maximum posterior probability (e.g. $P(w|x)$).
- Disadvantage: This model is a little less simplistic to the posterior probability model ($P(w|x)$).
- Compare/contrast: Part 1's explanation.

Part 3: $P(w|x) \geq P(x|w)P(w)/P(x)$ w.r.t. Generative, Discriminative, and Discriminate functions

- Generative models: use data available to estimate probabilities of each quantity: priors $P(w)$, likelihood $P(x|w)$, and evidence $P(x)$. This model is best for the model $P(x|w)P(w)/P(x) \Rightarrow$ joint distribution $P(x, w)$ and normalization to obtain posterior probabilities.
- Discriminative models: are based on the posterior probability $P(w|x)$ which do not attempt to fully model joint distributions $P(w, x)$, but attempts to calculate/estimate posterior directly, thus it is better for the $P(w|x)$ models.
 - \rightarrow Cannot do outlier detection like generative model
- Discriminate functions: Neither $P(w|x)$ or $P(x|w)P(w)/P(x)$ models this particularly well because these functions focus on decision boundaries which bypass likelihood probabilities & approaches.
 - \therefore For the $P(x|w)P(w)/P(x)$ model describing $P(x)$ can be beneficial because this tells us the probability of event x happening across all other cases (provides us with more beneficial data). $P(x)$ is said to describe the evidence.