

EECE-5644
Intro. to Machine Learning
SPRING 2016
Homework #: 2

Submission Date: 2 / 8 / 2016

Name (Last, First):

Homework Submission Rules:

1. All submitted work should be legible. Do not write in small or in overly cursive characters.
 2. Your submission of problem solutions must be in the given order, i.e., P1, P2, P3, etc. Do not submit in a random order.
 3. Use this cover page as the first page for each homework submission. Homework submitted without a cover page will result in a 10% reduction in the overall score for that homework.
 4. All plots must have axis labeled, titles, and legends if applicable.
 5. Code specified in the problems must be submitted via the Blackboard site alongside a pdf of your homework solutions. The code must be contained in a zip file with the following file name format: hw_##_lastname.zip. Do not print the code
 6. Be concise when writing your solutions and use both sides of the page.

(PLEASE DO NOT WRITE BELOW THIS LINE)

Elise Jortberg

EECE 5694

HW 2

1) 2.12

$$P(w_{\max}|x) \geq P(w_i|x) \text{ for } i=1\dots c$$

a) Show $P(w_{\max}|x) \geq \frac{1}{c}$

$$\int_0^\infty p(x) dx = 1 \rightarrow \sum_{i=1}^c p(w_i|x) = 1$$

upper limit $p(w_i|x) = p(w_{\max}|x)$

$$\therefore p(w_{\max}|x) \geq \frac{1}{c} \text{ else would never add to 1}$$

• If $p(w_{\max}|x) < \frac{1}{c}$ then $\sum = \frac{1}{c} + \frac{1}{c+a_1} + \frac{1}{c+a_2} + \dots$

never converges to 1

"worst case" scenario, all $p(w_i|x)$ equal

then $p(w_{\max}|x) = \frac{1}{c}$, so otherwise has to

$$be \geq \frac{1}{c}$$

b) Show $p(\text{error}) = 1 - \int p(w_{\max}|x)p(x) dx$

$$p(\text{error}) = 1 - \underbrace{p(\text{correct})}_{= \text{likelihood} \cdot \text{prior}}$$

$$= 1 - \int p(w_{\max}|x)p(x) dx$$

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

$$p(\text{error}) \leq 1 - \int \frac{1}{c} p(x) dx$$

$$\leq 1 - \frac{1}{c} \int p(x) dx,$$

$$< 1 - \frac{1}{c},$$

$$\leq \frac{c-1}{c}$$

d) The error in c) would occur when $p_1(x) = p_2(x)$ (distributions are the same). i.e. a gaussian where $\mu_1 = \mu_2$ and $\Sigma_1 = \Sigma_2$

2) 3.1

$$p(x|\theta) = \begin{cases} \theta e^{-\theta x} & x \geq 0 \\ 0 & \text{else} \end{cases}$$

a) see figures

b) Show $\hat{\theta}_{ML} = \frac{1}{n} \sum_{k=1}^n x_k$

$$\nabla_{\theta} p(D|\theta) = 0$$

$$D = \{x_1, \dots, x_n\}$$

$$\hat{\theta}_{ML} = \operatorname{argmax}_{\theta} \ln(p(D|\theta))$$

$$\sum_{k=1}^n \ln(p(x_k|\theta)) = \sum_{k=1}^n \ln(\theta) - \theta x_k$$

$$\nabla_{\theta} (n \ln \theta - \sum_{k=1}^n \theta x_k) = 0$$

$$\frac{\partial}{\partial \theta} (n \ln \theta - \sum_{k=1}^n \theta x_k) = 0$$

$$\frac{n}{\theta} - \sum_{k=1}^n x_k = 0 \rightarrow \frac{n}{\theta} = \sum_{k=1}^n x_k$$

$$\hat{\theta}_{ML} = \frac{n}{\sum_{k=1}^n x_k} = \frac{1}{n} \sum_{k=1}^n x_k$$

3) 3.2

$$p(x|\theta) \sim U(0, \theta) = \begin{cases} 1/\theta & 0 \leq x \leq \theta \\ 0 & \text{else} \end{cases}$$

a) $l(\theta) = \prod_{i=1}^n p(x_i|\theta)$

$$= \prod_{i=1}^n \frac{1}{\theta} \quad \text{where } x_i \leq \theta = \frac{1}{\theta^n}$$

$l(\theta)$

$$D = \{x_1, \dots, x_n\}$$

$$x_1 \ x_2 \ x_3 \ \dots \ x_n \ \theta$$

all x 's $< \theta$

$\theta < \max(D)$ when $l(\theta) = 0$
 $\theta \geq \max(D)$ when $l(\theta) = \frac{1}{\theta^n}$

$l(\theta) = \frac{1}{\theta^n}$ is maximized when θ is small
 $\therefore \hat{\theta}_{ML} = \max(D) = x_n$

4) 3.3

$$\begin{cases} z_{ik} = 1 & \text{k-th sample} = w_i \\ z_{ik} = 0 & \text{else} \end{cases}$$

a) Show $P(z_{i1}, \dots, z_{in} | p(w_i)) = \prod_{k=1}^n p(w_i)^{z_{ik}} (1-p(w_i))^{1-z_{ik}}$

$$\begin{cases} p(w_i) = p(z_{ik} = 1 | p(w_i)) & \text{if: } z_{ik} \\ 1-p(w_i) = p(z_{ik} = 0 | p(w_i)) & \text{else: } 1-z_{ik} \end{cases}$$

$p(z|w)$ = ? for all z_{ik}

in general: $p(x|\theta) = \prod_{k=1}^n p(x_k|\theta)$ (1)

$$\nabla_{\theta} \ln(p(x|\theta)) = 0 \quad (2)$$

$$p(z_{ik} | p(w_i)) = p(w_i)^{z_{ik}} (1-p(w_i))^{1-z_{ik}}$$

$$\therefore P(z_{i1}, \dots, z_{in} | p(w_i)) = \prod_{k=1}^n p(w_i)^{z_{ik}} (1-p(w_i))^{1-z_{ik}}$$

b) Show $\hat{P}_{mc}(w_i) = \frac{1}{n} \sum_{k=1}^n z_{ik}$

$$(2): \frac{\partial}{\partial p(w_i)} \ln \left(\prod_{k=1}^n p(w_i)^{z_{ik}} (1-p(w_i))^{1-z_{ik}} \right)$$

$$\frac{\partial}{\partial p(w_i)} \sum_{k=1}^n z_{ik} \ln(p(w_i)) + (1-z_{ik}) \ln(1-p(w_i)) = 0$$

$$\sum_{k=1}^n z_{ik} \frac{1}{p(w_i)} + \frac{1}{1-p(w_i)} \sum_{k=1}^n (1-z_{ik}) = 0$$

$$\sum_{k=1}^n z_{ik} \frac{1}{p(w_i)} = \frac{1}{1-p(w_i)} \sum_{k=1}^n (1-z_{ik})$$

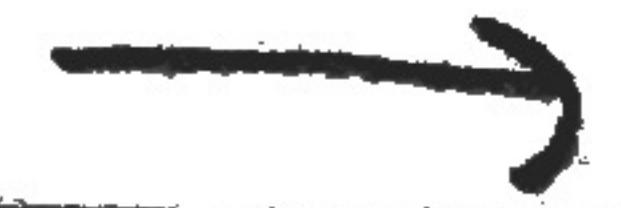
$$(1-p(w_i)) \sum_{k=1}^n z_{ik} = p(w_i) \sum_{k=1}^n (1-z_{ik})$$

$$\sum_{k=1}^n z_{ik} - p(w_i) \sum_{k=1}^n z_{ik} = p(w_i) \sum_{k=1}^n 1 - p(w_i) \sum_{k=1}^n z_{ik}$$

$$\sum_{k=1}^n z_{ik} - p(w_i) \sum_{k=1}^n z_{ik} = n p(w_i) - p(w_i) \sum_{k=1}^n z_{ik}$$

$$\sum_{k=1}^n z_{ik} (1-p(w_i)) = p(w_i) (n - \sum_{k=1}^n z_{ik})$$

$$\frac{\sum_{k=1}^n z_{ik}}{n - \sum_{k=1}^n z_{ik}} = \frac{p(w_i)}{1-p(w_i)}$$



$$\therefore p(w_i) = \frac{1}{n} \sum_{k=1}^n z_{ik}$$

This result means that we estimate/guess at $p(w_i)$ based on the input/training dataset to initialize the classifier/model. We don't base guess on prior/human knowledge.

5) 3.4

x = binary vector w/ distribution

$$P(x|\vec{\theta}) = \prod_{i=1}^d \theta_i^{x_i} (1-\theta_i)^{1-x_i}$$

where $\vec{\theta} = (\theta_1, \dots, \theta_d)^T$ unknown parameter vector

θ_i = probability $x_i = 1$

$$\text{Show } \hat{\vec{\theta}}_{\text{ML}} = \frac{1}{n} \sum_{k=1}^n \vec{x}_k$$

in general: $p(D|\theta) = \int p(x|\theta) p(\theta|D) d\theta$

$$D = \{x_1, \dots, x_n\} \rightarrow \text{multivariate}$$

$$\theta = \{\theta_1, \dots, \theta_d\}$$

$$p(D|\theta) = \prod_{k=1}^n p(x_k|\theta) \text{ becomes}$$

$$p(D|\theta) = \prod_{k=1}^n \prod_{i=1}^d \theta_i^{x_{ki}} (1-\theta_i)^{1-x_{ki}}$$

$$\nabla_{\theta} \ln(p(D|\theta)) = 0$$

$$\nabla_{\theta} \cdot \sum_{k=1}^n \sum_{i=1}^d x_{ki} \ln(\theta_i) + (1-x_{ki}) \ln(1-\theta_i) = 0$$

$$\sum_{k=1}^n \sum_{i=1}^d \frac{x_{ki}}{\theta_i} + (1-x_{ki}) \frac{1}{1-\theta_i}$$

find each $\frac{\partial \ln p(x|\theta)}{\partial \theta_j}$ based on arbitrary j for singular case

$$\sum_{j=1}^n \sum_{i=1}^d \frac{\partial}{\partial \theta_i} [\bar{x}_{ji} \ln \theta_i + (1-\bar{x}_{ji}) \ln(1-\theta_i)]$$

$$\frac{\partial}{\partial \theta_i} = \bar{x}_{j2} \ln \theta_2 = 0$$

$$\frac{\partial}{\partial \theta_i} = \bar{x}_{j1} \ln \theta_1 = \frac{\bar{x}_{j1}}{\theta_1}$$

substitute in so far one case
other thetas $\rightarrow 0$

\therefore for any i in vector $\theta = \{\theta_1, \dots, \theta_d\}$

$$[\nabla_{\theta} l(\theta)]_i = \nabla_{\theta_i} l(\theta)$$

$$\frac{1}{\theta_i} \sum_{k=1}^n x_{ki} = \frac{1}{1-\theta_i} \sum_{k=1}^n (1-x_{ki})$$

Solve for $\hat{\theta}_i$:

$$1 - \hat{\theta}_i \sum_{k=1}^n x_{ki} = \hat{\theta}_i \sum_{k=1}^n (1-x_{ki}) \quad (\text{like problem 3.3})$$

$$\therefore \hat{\theta}_i = \frac{1}{n} \sum_{k=1}^n x_{ki}$$

since we solved for arbitrary i , this applies
for all $i \in [1, d]$

$$\vec{\theta}_m = \frac{1}{n} \sum_{k=1}^n \vec{x}_k$$

(e) 3.17

Derive Bayesian classifier for d dimensional multivariate Bernoulli case. $P(\vec{x} | D) = P(\vec{x} | \vec{\theta}, w_i)$

$$P(\vec{x} | \vec{\theta}) = \prod_{i=1}^d \theta_i^{x_i} (1-\theta_i)^{1-x_i}$$

$$D = \{\vec{x}_1, \dots, \vec{x}_n\}$$

a) if $\vec{s} = (s_1, \dots, s_d)^T$ is sum of n samples show

$$P(D | \vec{\theta}) = \prod_{i=1}^d \theta_i^{s_i} (1-\theta_i)^{n-s_i}$$

$$17) p(\vec{x}|\vec{\theta}) = \prod_{i=1}^d \theta_i^{x_i} (1-\theta_i)^{1-x_i}$$

$$D = \{\vec{x}_1, \dots, \vec{x}_n\}$$

a) $\vec{s} = (s_1, \dots, s_d)^T$ sum of n samples

$$p(D|\vec{\theta}) = \prod_{k=1}^n p(\vec{x}_k|\vec{\theta})$$

k has associated d , k is independent

$$p(\vec{x}_k|\vec{\theta}) = \prod_{i=1}^d \theta_i^{x_{ki}} (1-\theta_i)^{1-x_{ki}}$$

$$p(D|\vec{\theta}) = \prod_{k=1}^n \prod_{i=1}^d \theta_i^{x_{ki}} (1-\theta_i)^{1-x_{ki}}$$

for $\theta^{x_{ki}}, (1-\theta_i)^{1-x_{ki}}$ multiplying over $k, 1 \rightarrow n$

$$p(D|\vec{\theta}) = \prod_{i=1}^d \theta^{\sum_{k=1}^n x_{ki}} (1-\theta)^{\sum_{k=1}^n (1-x_{ki})}$$

$$p(D|\vec{\theta}) = \prod_{i=1}^d \theta^{s_i} (1-\theta)^{n-s_i}$$

b) assume uniform prior distribution for $\vec{\theta}$ and
identity $\int_0^1 \theta^m (1-\theta)^n d\theta = \frac{m! n!}{(m+n+1)!}$ solve $p(\vec{\theta}|D)$

$$p(\vec{\theta}|D) = \frac{p(D|\vec{\theta}) p(\vec{\theta})}{p(D)}$$

uniform $\therefore p(\vec{\theta}) = 1 \rightarrow$ need $p(D) \dots$

$$p(D) = \int p(D|\vec{\theta}) p(\vec{\theta}) d\vec{\theta}$$

$$= \int_{\text{to } d} \prod_{i=1}^d \theta^{s_i} (1-\theta)^{n-s_i} d\theta_{(1, \dots, d)}$$

multiple integrals $\rightarrow d = \int_{\text{to } d}, \text{ BINARY}$

$$p(D) = \prod_{i=1}^d \int_0^1 \theta^{s_i} (1-\theta)^{n-s_i} d\theta_i \quad \text{for arbitrary case}$$

identity: $m = s_i, n = n - s_i$

$$p(D) = \prod_{i=1}^d \frac{s_i! (n-s_i)!}{(s_i + n - s_i + 1)!} = \prod_{i=1}^d \frac{s_i! (n-s_i)!}{(n+1)!}$$

plug back in to $p(\vec{\theta}|D)$ to solve \rightarrow

$$p(\vec{\theta} | D) = \frac{\prod_{i=1}^d \theta_i^{s_i} (1-\theta_i)^{n-s_i}}{\frac{\prod_{i=1}^d s_i! (n-s_i)!}{(n+1)!}}$$

$$= \prod_{i=1}^d \frac{(n+1)! \theta_i^{s_i} (1-\theta_i)^{n-s_i}}{s_i! (n-s_i)!}$$

c) plot $p(\vec{\theta} | D)$ for $d=1, n=1$ and both cases of s

$$p(\vec{\theta}_{n=1} | D) = \frac{2!}{s_1! (1-s_1)!} \theta_1^{s_1} (1-\theta_1)^{1-s_1}$$

$$s_1 = \begin{cases} 0 & \text{for } 0 \leq \theta_1 \leq 1 \\ 1 & \end{cases}$$

FIGURE IN MATLAB

d) Integrate $p(\vec{x} | \vec{\theta}) p(\vec{\theta} | D) d\theta$ to find $p(\vec{x} | D)$

$$p(\vec{x} | D) = \int p(\vec{x} | \vec{\theta}) p(\vec{\theta} | D) d\theta$$

$$= \int \prod_{i=1}^d \theta_i^{x_{ki}} (1-\theta_i)^{1-x_{ki}} \prod_{i=1}^d \frac{(n+1)! \theta_i^{s_i} (1-\theta_i)^{n-s_i}}{s_i! (n-s_i)!} d\theta$$

$$= \int \prod_{i=1}^d \theta_i^{x_i+s_i} (1-\theta_i)^{1-x_i+n-s_i} \underbrace{\frac{(n+1)!}{s_i! (n-s_i)!}}_{\text{no } \theta \text{ component}} d\theta$$

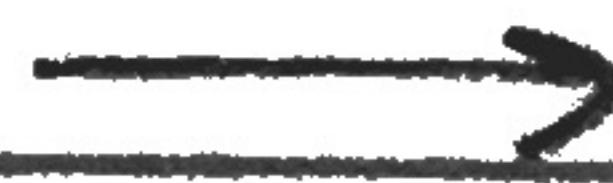
$$= \prod_{i=1}^d \frac{(n+1)!}{s_i! (n-s_i)!} \int_0^1 \theta_i^{x_i+s_i} (1-\theta_i)^{1-x_i+n-s_i} d\theta$$

$$= \prod_{i=1}^d \frac{(n+1)!}{s_i! (n-s_i)!} \frac{(x_i+s_i)! (1-x_i+n-s_i)!}{(x_i+s_i + 1 - x_i + n - s_i + 1)}$$

$$p(\vec{x} | D) = \prod_{i=1}^d \frac{(n+1)!}{s_i! (n-s_i)!} \frac{(x_i+s_i)! (1-x_i+n-s_i)!}{(n+2)!}$$

$$\frac{(n+1)!}{(n+2)!} = \frac{1}{n+2}$$

$$p(\vec{x} | D) = \frac{(x_i+s_i)! (1-x_i+n-s_i)!}{s_i! (n-s_i)! (n+2)}$$



if $x_i = 1$

$$p(x|D) = \frac{(1+s_i)! (1-n-s_i)!}{s_i! (n-s_i)! (n+2)} = \frac{s_i + 1}{n+2}$$

if $x_i = 0$, i.e. ($x_i = 1$)

$$p(x|D) = 1 - \frac{s_i + 1}{n+2}$$

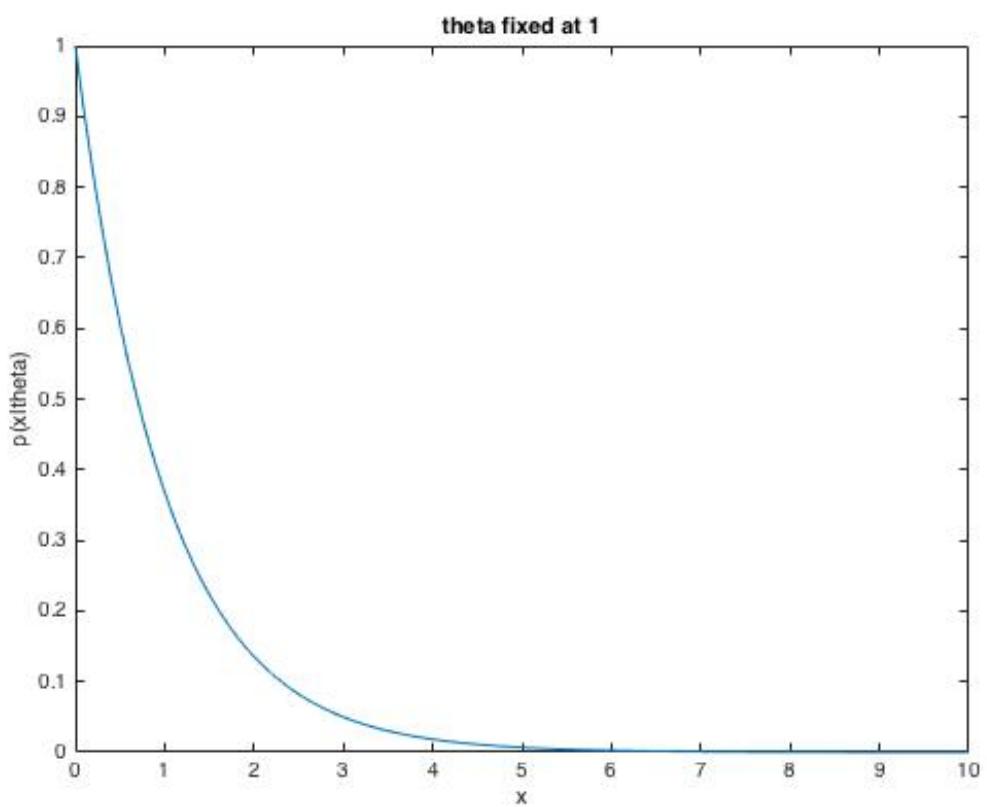
$$p(\vec{x}|D) = \left(\frac{s_i + 1}{n+2}\right)^{x_i} \left(1 - \frac{s_i + 1}{n+2}\right)^{1-x_i}$$

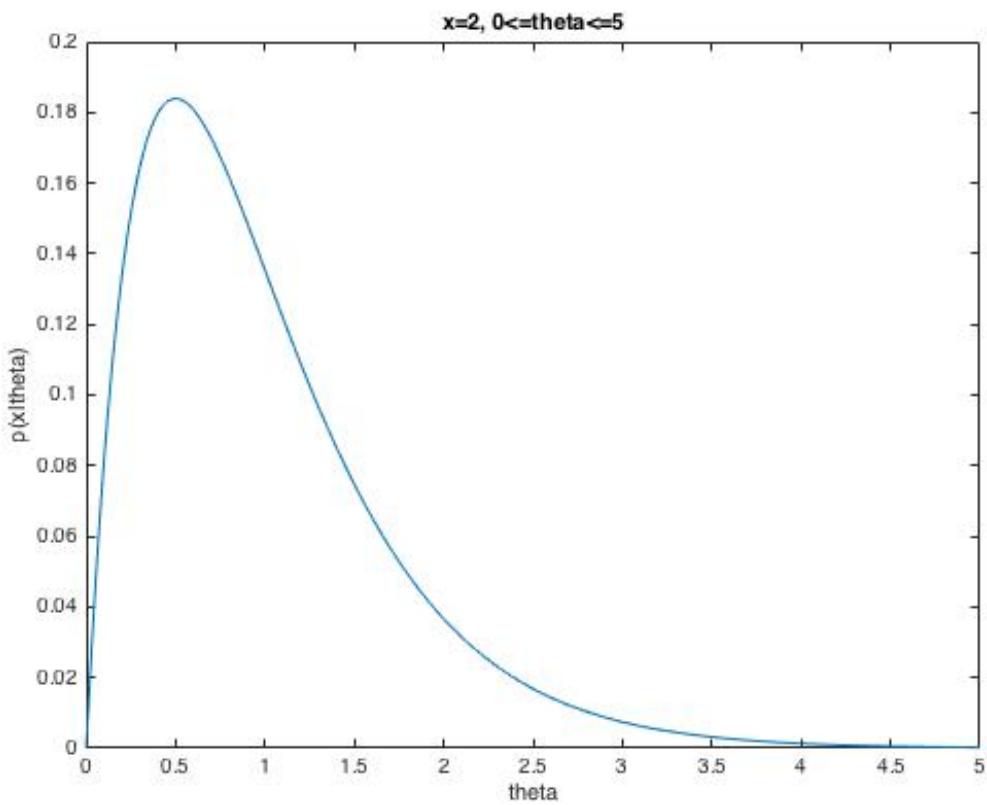
e) obtain $p(\vec{x}|D)$ by using $\hat{\theta}$ for $\vec{\theta}$ in $p(\vec{x}|\vec{\theta})$
what is effective Bayesian estimate for $\vec{\theta}$?

$p(\vec{x}|D)$ boils down to $\hat{\theta}_i = \frac{s_i + 1}{n+2}$ to

Predict $\hat{\theta}_i$ for x_i case (not using $(1-\hat{\theta}_i)^{1-x_{ki}}$ term).

Problem 3.1:





Problem 3.17

