

Midterm exam for Fundamentals of Statistical Machine Learning (201910). Open book test.

Your Name: \_\_\_\_\_

1. Consider a computer program C which outputs 'a' or 'b' each time it runs. It is known that C uses a random generator to output a with probability  $p_a$ .

No one knows the value of  $p_a$ , but you have a friend, Thomas, who plans to run the program some number of times and then predict the output of the next run of the program.

Thomas loves Bayesian statistics, so his posterior probability estimate certainly will include a prior.

Unfortunately Thomas is coy and will not tell you what his prior is.

Instead he tells you that if he were to run the program 4 times and get an a each time, he would be 90% sure that next run would also output an a. Likewise, if he were to get 4 out of 4 b's, he would be 90% sure that the next run would output another b.

**Question:** Is Thomas using a uniform prior? Explain your answer.

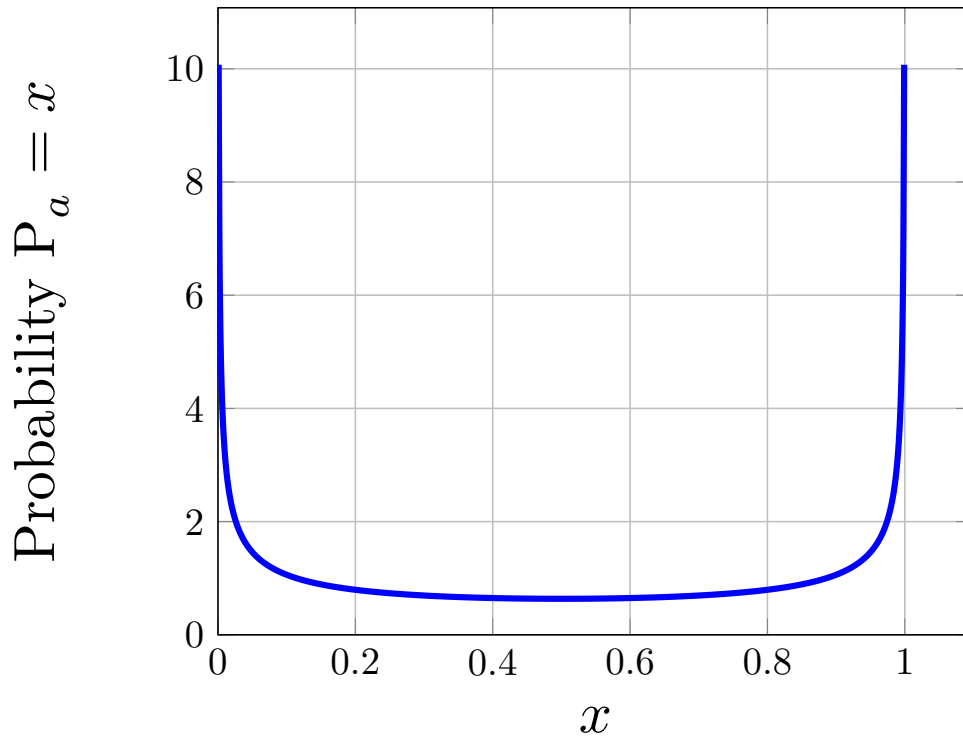


Figure 1: Beta( $1/2, 1/2$ ) — Jeffrey's Prior

**Solution:**

Thomas is using a Jeffreys Prior, which is not the same as a uniform prior. A uniform prior stipulates a pseudocount of one for each class (in this case **a** or **b**), so after 4 **a**'s in a row, the posterior estimate of the probability of getting another **a** would be  $\frac{4+1}{4+2} = 5/6 \neq 90\%$ .

Instead Thomas is using a pseudocount of  $1/2$  for both **a** and **b**, yielding  $\frac{4^{1/2}}{4+1} = 90\%$ . This is the Jeffreys prior.

$$\text{Beta}(x|a, b) \stackrel{\text{def}}{=} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$$

A uniform prior is equivalent to  $\text{Beta}(1, 1) = 1x^{1-1}(1-x)^{1-1} = 1$ .

The Jeffreys prior is  $\text{Beta}(1/2, 1/2)$  with pdf:

$$\text{Beta}(1/2, 1/2) = \frac{\Gamma(1/2 + 1/2)}{\Gamma(1/2)\Gamma(1/2)} x^{1/2-1} (1-x)^{1/2-1} = \frac{1}{\sqrt{(\pi)}\sqrt{(\pi)}} x^{-1/2} (1-x)^{-1/2} = \frac{1}{\pi\sqrt{x(1-x)}}$$

Your Name: \_\_\_\_\_

2. This question asks about the entropy of a probability distribution in which two categories have been joined together.

Let  $F$  denote a probability distribution over the first  $k$  natural numbers:  $1, 2, \dots, k$ . Denote the probability  $F(i)$  as  $p_i$ , so  $\sum_1^k p_i \equiv 1$ .

(i.e.  $F$  is like a die (骰子) with  $k$ -sides, each side having its own probability  $p_i$ )

Let  $F'$  be a probability distribution over  $1, 2, \dots, k-1$ ;

almost the same as  $F$ , but with the last two elements ( $k-1$  and  $k$ ) merged.

So,  $p'_i = p_i$   $1 \leq i < k-1$ , and  $p'_{k-1} = p_{k-1} + p_k$ . There is no  $p'_k$  (or equivalently  $p'_k \stackrel{\text{def}}{=} 0$ ).

Let  $H()$  denote information theoretic entropy.

GIVEN: We know the theoretic entropy of  $F$ ,  $H(F) = 5.0$ ; and that  $p_{k-1} = p_k = 0.01$ .

1. Give a lower bound on  $k$  and explain why  $k$  must be at least that large.
2. What is  $H(F')$ , the entropy of  $F'$ ?

**Solution:**

1. An entropy of  $h$  bits is the amount of information needed to select among  $2^h$  equally likely things. So one way to have  $H(F) = 5$  is the uniform distribution over  $1, 2, \dots, 32$ . Therefore  $k$  could be as small as 32. Moreover for a fixed  $k$ , a uniform distribution maximizes entropy, so for example if  $k$  were 31, the entropy would have to be less than  $\lg 31 \approx 4.95 < 5$ . Thus 32 is a lower bound on  $k$ .

2. Using the “increases incrementally when splitting categories” property of entropy, We have:

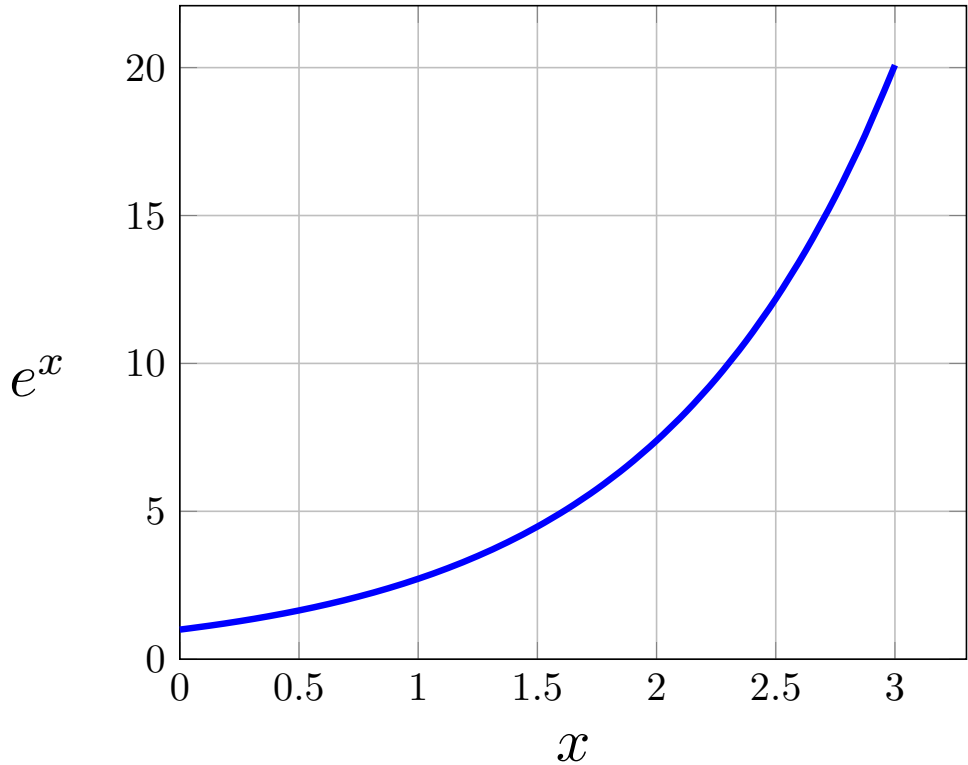
$$H(F) = H(F') + p'_{k-1} H(0.5, 0.5) \tag{1}$$

Rearranging,

$$H(F') = H(F) - p'_{k-1} H(0.5, 0.5) = 5 - 0.02 \cdot 1 = 4.98 \tag{2}$$

Because  $H(0.5, 0.5)$  is a 50%-50% split, which is one bit of information.

Your Name: \_\_\_\_\_



3.

Consider a mixture model of two normal (one-dimensional Gaussian) distributions:

$$M(x) = \frac{2}{3}N_a + \frac{1}{3}N_b$$

Denoted  $N_a : N(\mu_a, \sigma_a^2)$  and  $N_b : N(\mu_b, \sigma_b^2)$ .

In other words, when generating points from M:

first one of the two components  $N_a$  or  $N_b$  are selected are random with  $N_a$  selected  $\frac{2}{3}$  of the time, and then a point is generated according to the mean and variance of the selected component.

A data point  $x_1$  is sampled from  $M(x)$ .

Consider the probability that  $x_1$  "came from"  $N_a$ .

1. Give the general formula for that probability (a formula including  $x_1, \dots$ )
2. Give the approximate numerical value for the probability for the case:

$$\begin{array}{ccccc} x_1 & \mu_a & \sigma_a & \mu_b & \sigma_b \\ 4 & -1 & 2 & 18 & 4 \end{array}$$

Approximation should be made *without* using a calculator (不可以用計算機) and should be within  $\pm 10\%$ .

The answer can be given as a fraction, e.g.  $\frac{2}{7}$  instead of 0.285714....

**Solution:**

For convenience, first define  $z_a = (\mu_a - x_1)/\sigma_a$ ,  $z_b = (\mu_b - x_1)/\sigma_b$ ;  
and compute the probability ratio  $\frac{P[N_a|x]}{P[N_b|x]}$  first.

**General Formula**

$$\frac{P[N_a|x]}{P[N_b|x]} = \frac{P[N_a]^{1/\sigma_a} \text{EXP}[-z_a^2/2]}{P[N_b]^{1/\sigma_b} \text{EXP}[-z_b^2/2]} = \frac{P[N_a] \sigma_b \text{EXP}[z_b^2/2]}{P[N_b] \sigma_a \text{EXP}[z_a^2/2]} = \frac{P[N_a] \sigma_b}{P[N_b] \sigma_a} \text{EXP}\left[\frac{z_b^2 - z_a^2}{2}\right]$$

Finally convert the ratio into a probability via:

$$P[A] \equiv \frac{1}{1 + \frac{P[B]}{P[A]}} \text{ obtaining } P[N_a|x] = \frac{1}{1 + \frac{P[N_b] \sigma_a}{P[N_a] \sigma_b} \text{EXP}\left[\frac{z_a^2 - z_b^2}{2}\right]}$$

**Numerical Answer** Adding to the table in the question:

$x_1$	$\mu_a$	$\sigma_a$	$z_a$	$\mu_b$	$\sigma_b$	$z_b$	$\frac{P[N_a]}{P[N_b]}$	$\frac{z_b^2 - z_a^2}{2}$	$\text{EXP}\left[\frac{z_b^2 - z_a^2}{2}\right]$
4	-1	2	5/2	18	4	-7/2	2	$\frac{(-7)^2 - (5)^2}{2} = \frac{49 - 25}{2} = \frac{24}{2} = 3$	$\text{EXP}[3] \approx 20$

$$\frac{P[N_a|x]}{P[N_b|x]} = \frac{P[N_a] \sigma_b}{P[N_b] \sigma_a} \text{EXP}\left[\frac{z_b^2 - z_a^2}{2}\right] \approx 2 \cdot \frac{4}{2} \cdot 20 \approx 80 \text{ so } P[N_a] \approx \frac{80}{81}$$

Your Name: \_\_\_\_\_

Consider a typically classification problem. For example one with  $m$  features  $F_1, \dots, F_m$ , which may be symptoms, e.g. fever, coughing, etc.; while the class to be predicted is the person's condition (healthy or disease, etc.).

Consider the following statements:

$$P[F_1, \dots, F_m] \Leftrightarrow \prod_{i=1}^m P[F_i] \tag{1}$$

$$P[F_1, \dots, F_m|C] \Leftrightarrow P[F_1|C] P[F_2|C, F_1] P[F_3|C, F_1, F_2] \dots P[F_m|C, F_1, \dots, F_{m-1}] \tag{2}$$

$$P[C|F_1, \dots, F_m] \Leftrightarrow P[C, F_1] P[C, F_2, F_1] P[C, F_3, F_1, F_2] \dots P[C, F_m, F_1, \dots, F_{m-1}] \tag{3}$$

$$P[C|F_1, \dots, F_m] \Leftrightarrow \prod_{i=1}^m P[C, F_i] \tag{4}$$

$$P[F_1, \dots, F_m|C] \Leftrightarrow \prod_{i=1}^m P[F_i|C] \tag{5}$$

For each blank in the table below fill in one of {always, often, seldom, hardly}, where **always** means always equal, **often** means often approximately equal, **seldom** means usually unequal, and **hardly** means hardly ever equal.

Eq num:	(1)	(2)	(3)	(4)	(5)
Equal?					

**Why?** Give a the reason for your answers (可以用中文).

<b>Solution: Solution:</b>					
Eq num:	(1)	(2)	(3)	(4)	(5)
Equal?	seldom	always	hardly	hardly	often
Reasons					
<p>(1) Assumes features are independent of each other. This is possible, but <b>seldom</b> happens because in a typical classification problem the features tend to depend on each other via the class. For example people with fevers are more likely to cough, because those two features both correlate with the health of the person.</p> <p>(2) Is an instance of the “chain rule” for conditional probability and is identically true.</p> <p>(3) Is hardly ever true. Note that in general the right hand side of equation (3) gets smaller and smaller with increasing <math>m</math>, but the left hand side doesn't.</p> <p>(4) Is also hardly ever true. Same reason as for (3).</p> <p>(5) Is often approximately true. Note that this is the approximation behind the Naïve Bayes classifier, which often works reasonably well as a classifier.</p>					

Your Name: \_\_\_\_\_

	start	unknowns	finish
example 1.	$(a-b)(a+b)$	$\equiv \underline{?}^2 - \underline{?}^2$	$\Rightarrow a^2 - b^2$
example 2.	$x!$	$\equiv \Gamma(\underline{?})$	$\Rightarrow \Gamma(x+1)$
problem 1.	$\binom{n}{k}$	$\equiv \frac{\Gamma(\underline{?})}{\Gamma(\underline{?})\Gamma(\underline{?})}$	$\Rightarrow$ _____
problem 2.	Binomal( $k p, n$ )	$\underline{?}$ Beta( $\underline{?}, \underline{?}, \underline{?}$ )	$\Rightarrow$ _____ Beta(_____)

**Problem 1.** Fill in problem 1. in the table above relating  $\binom{n}{k}$  to  $\frac{\Gamma(\underline{?})}{\Gamma(\underline{?})\Gamma(\underline{?})}$ .

The Beta distribution is defined as:

$$\text{Beta}(x|a, b) \stackrel{\text{def}}{=} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$$

**Problem 2.** Express the Beta distribution parameters ( $a, b$ ) and data point  $x$  in terms of ( $k, p$  and  $n$ ), such that

$$\text{Binomal}(k|p, n) = F(k, p, n) \text{Beta}(x|a, b)$$

Where  $F(\cdot, p, n)$  is a simple function of  $k, p$ , and  $n$ .



**Solution:**

The Binomial and Beta distribution are defined as:

$$\text{Binomial}(k|p, n) \stackrel{\text{def}}{=} \binom{n}{k} p^k (1-p)^{n-k} \quad \text{Beta}(x|a, b) \stackrel{\text{def}}{=} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$$

We can see that  $x$  corresponds to  $p$ , so first rename  $p$  as  $x$  in Binomial(  $k$  |  $x$ ,  $n$  ).  
Also consider the normalization terms:

By inspecting the form of these equations, we see likely correspondences:  
 $p \Leftrightarrow x$ ,  $a \Leftrightarrow k$ ,  $b \Leftrightarrow n - k$ .

Trying this we obtain:

$$\begin{aligned} \text{Binomial}(k|p, n) &= \binom{n}{k} p^k (1-p)^{n-k} \\ &= \binom{a+b}{a} p^a (1-p)^b && \text{Substituting } k \rightarrow a, n-k \rightarrow b: n \rightarrow a+b \\ &= \frac{(a+b) \Gamma(a+b)}{ab \Gamma(a) \Gamma(b)} p^{a-1} p (1-p)^{b-1} (1-p) \\ &= \frac{(a+b) \Gamma(a+b)}{ab \Gamma(a) \Gamma(b)} p^{a-1} (1-p)^{b-1} p(1-p) \\ &= \frac{(a+b) \text{Beta}(p|a, b) p(1-p)}{ab} && \text{Definition of Beta}(p|a, b) \\ &= p(1-p) \frac{n}{k(n-k)} \text{Beta}(p|k, n-k) && \text{Restoring: } a \rightarrow k, b \rightarrow n-k \text{ and rearranging} \end{aligned}$$