

## Homework 4 Solutions

### 1. Using conditional probability for estimation

Suppose  $x : \Omega \rightarrow U$  and  $w : \Omega \rightarrow V$  are independent random variables with pmfs  $p^x$  and  $p^w$ . Define the random variable  $y$  by

$$y = x + w$$

We would like to measure  $y$  and hence determine  $x$ , without measuring  $w$ . We have

$$V = \{1, 2, 3, 4, 5\} \quad \text{and} \quad U = \{1, 2, 3, 4, 5, 6\}$$

and

$$p^w = \begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.2 \\ 0.2 \end{bmatrix} \quad \text{and} \quad p^x = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.15 \\ 0.15 \\ 0.25 \\ 0.25 \end{bmatrix}$$

(a) Show that

$$\text{Prob}(x = a \text{ and } y = b) = p^x(a)p^w(b - a)$$

where  $p^w(b - a)$  is understood to be zero for  $b - a < 1$  or  $b - a > 5$ .

(b) Find the induced pmf of  $y$ .

(c) Now suppose we measure  $y_{\text{meas}} = 7$ . Find the conditional pmf  $q$  where

$$q(a) = \text{Prob}(x = a \mid y = y_{\text{meas}})$$

Plot  $q$  against  $U$ .

(d) Now simulate the random variables  $w$  and  $x$ , and perform  $m = 10000$  trials. Keep data for  $x$  and  $w$ , and for each trial, compute  $y = w + x$ .

Compute the observed frequencies of  $y$ , and plot these. Compare with the pmf for  $y$ .

(e) Now look at the same data set, and write a loop that discards those data items for which  $y \neq y_{\text{meas}}$ . On the remaining data, compute the observed frequencies of  $x$ . Plot the observed frequencies of  $x$ , and compare with the conditional pmf of  $x$  given  $y_{\text{meas}}$ .

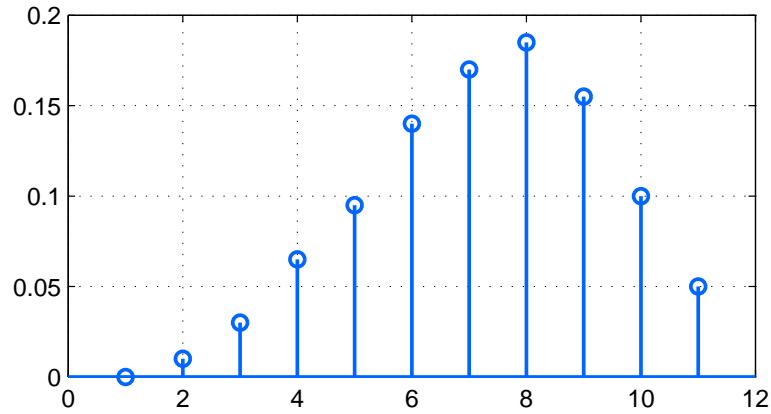
(f) Now you have to pick one estimate of  $x$  based on your measurement. Which is the estimate that minimizes your probability of error?

### **Solution.**

(a) We have

$$\begin{aligned} \text{Prob}(x = a \text{ and } y = b) &= \text{Prob}(y = b \mid x = a) \text{Prob}(x = a) \\ &= \text{Prob}(x + w = b \mid x = a) \text{Prob}(x = a) \\ &= \text{Prob}(w = b - a \mid x = a) \text{Prob}(x = a) \\ &= \text{Prob}(w = b - a) \text{Prob}(x = a) \quad \text{since } x \text{ and } w \text{ are independent} \\ &= p^x(a)p^w(b - a) \end{aligned}$$

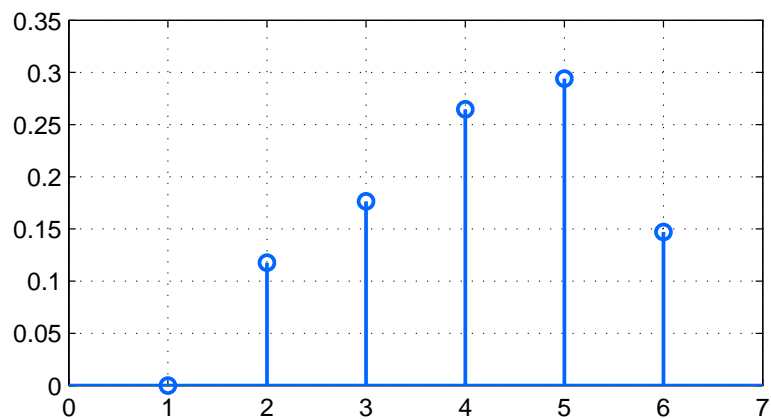
(b) We have  $y : \Omega \rightarrow W$  where  $W = [2, 3, \dots, 11]$ . The induced pmf of  $y$  is given by the convolution  $p^y = p^x * p^w$ . It is plotted below.



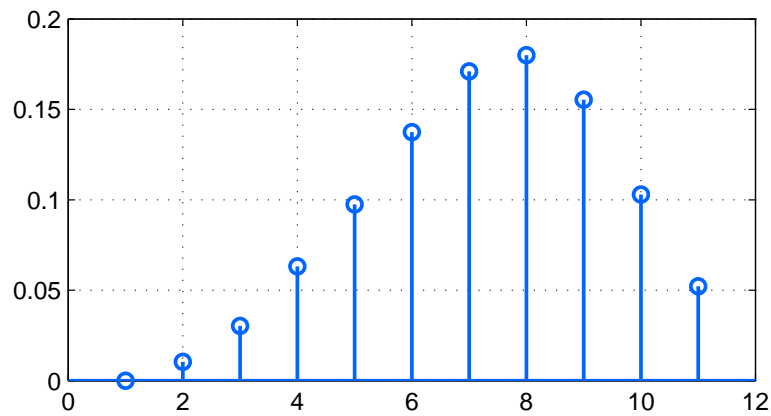
(c) The conditional pmf is given by

$$\begin{aligned}
 q(a) &= \text{Prob}(x = a \mid y = y_{\text{meas}}) \\
 &= \frac{\text{Prob}(x = a \text{ and } y = y_{\text{meas}})}{\text{Prob}(y = y_{\text{meas}})} \\
 &= \frac{p^x(a)p^w(y_{\text{meas}} - a)}{p^y(y_{\text{meas}})}
 \end{aligned}$$

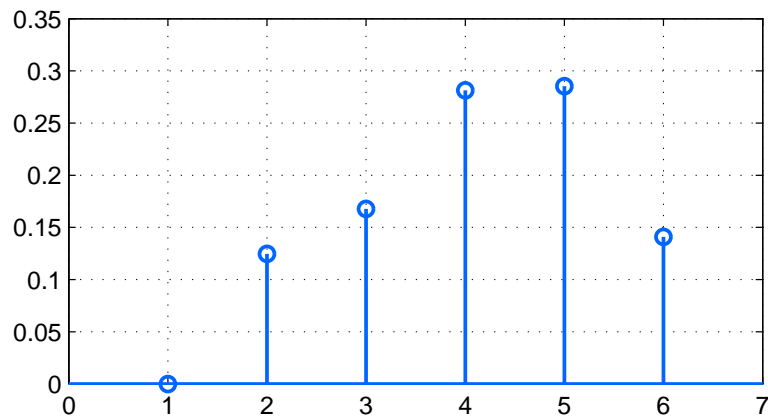
It is plotted below.



(d) The frequencies for  $y$  are shown below; they match the pmf for  $y$  as expected.



- (e) The plot of the observed frequencies of  $x$  when  $y = y_{\text{meas}}$  is shown below. It also matches the pmf  $q$  very well.



- (f) If  $y = y_{\text{meas}}$ ,  $\hat{x} = 5$  minimizes the probability of error, since it has the largest a-posteriori conditional probability.

## 2. Radar

A radar has conditional probability mass functions

$$p^{\text{nothing}} = 10^{-2} [4 \ 7 \ 11 \ 12 \ 17 \ 16 \ 13 \ 9 \ 6 \ 4 \ 1 \ 0 \ 0]$$

$$p^{\text{something}} = 10^{-2} [0 \ 0 \ 2 \ 6 \ 8 \ 10 \ 15 \ 15 \ 15 \ 13 \ 9 \ 5 \ 2]$$

Here  $X_1$  is the event that nothing is visible, and  $X_2$  is the event that something is visible, and  $Y_i$  is the event that the radar returns  $i$  pings (at least one ping is always returned.) Then

$$\text{Prob}(Y_i | X_1) = p_i^{\text{nothing}} \quad \text{Prob}(Y_i | X_2) = p_i^{\text{something}}$$

The prior probabilities are

$$\text{Prob}(X_1) = 0.6 \quad \text{Prob}(X_2) = 0.4;$$

- (a) Find the classifier with the minimum probability of error, and the corresponding probability of a correct classification.
- (b) Suppose we significantly prefer false positives to false negatives, and so decided to minimize

$$J_1 + \mu J_2$$

where

$$J_1 = \text{Prob}(\text{false positive}) \quad J_2 = \text{Prob}(\text{false negative})$$

First, let's arbitrarily choose  $\mu = 3$ .

Find the optimal decision rule to minimize the weighted sum cost. For this rule, find the probability of a false positive error, a false negative error, and the probability of a correct result.

- (c) To do this more systematically, we need to plot the *operating characteristic*. As usual when plotting trade-off curves, it's a good idea to choose a wide range of values for  $\mu$ . Pick a suitable set of values for  $\mu$ , and for each compute the classifier  $K$  which minimizes  $J_1 + \mu J_2$ , as well as the optimal values of  $J_1$  and  $J_2$ . Plot  $J_1$  against  $J_2$ . As you can see, it's not a very good radar.
- (d) Now we would like to find the classifier which minimizes the probability of false positives, subject to the constraint that the probability of false negatives is less than 0.1. First find the smallest  $\mu$  you can use that achieves this level of probability of false negatives, then find the corresponding optimal classifier. Again, for this rule, find the probability of a false positive error, a false negative error, and the probability of a correct result.

- (e) In fact there are only a finite number of possible classifiers? How many are there? (Don't only count the threshold rules, count all of them.)
- (f) Write a routine that for each possible decision rule, computes the corresponding probability of false positives and false negatives. Plot another copy of the operating characteristic, and add to it all of these new points.
- (g) The plot you obtained in part (f) has a symmetry. State precisely what this symmetry is. Can you explain how it arises?
- (h) What is the *worst* possible classifier you could use; i.e., the one that achieves the maximum possible probability of error?
- (i) Now among *all* possible decision rules, find the one that achieves the smallest possible probability of false positives, subject to the constraint that the probability of false negatives is less than 0.1. What is the decision rule, the probability of a false positive error, a false negative error, and the probability of a correct result. Can you explain why this decision rule was not found by the approach in part (d).
- (j) Let's also explore what would have happened if the radar were better, i.e., the conditional pmfs were further apart. Suppose

$$\text{Prob}(Y_i | X_2) = p_{i-k}^{\text{something}}$$

i.e, the conditional probabilities given something is visible are shifted to the right by  $k$  positions. For each  $k \in \{0, 1, 2, \dots, 5\}$  plot the operating characteristic, (all on one plot).

**Solution.**

- (a) Here

$$C = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

The optimal estimator is a threshold rule; estimate  $X_1$ , if less than or equal to 7 pings are returned; else estimate  $X_2$ . The corresponding probability of error is 0.284.

- (b) Here

$$C = \begin{bmatrix} 0 & \mu \\ 1 & 0 \end{bmatrix}$$

When  $\mu = 3$ , we have

$$JC^T = \begin{bmatrix} 0.000 & 0.024 \\ 0.000 & 0.042 \\ 0.024 & 0.066 \\ 0.072 & 0.072 \\ 0.096 & 0.102 \\ 0.120 & 0.096 \\ 0.180 & 0.078 \\ 0.180 & 0.054 \\ 0.180 & 0.036 \\ 0.156 & 0.024 \\ 0.108 & 0.006 \\ 0.060 & 0.000 \\ 0.024 & 0.000 \end{bmatrix}$$

There are two possible classifiers which are optimal in this case.

The first classifier estimates  $X_1$  if 1, 2, 3, 5 pings are returned, else it estimates  $X_2$ . The error matrix is

$$E = \begin{bmatrix} 0.234 & 0.040 \\ 0.366 & 0.360 \end{bmatrix}$$

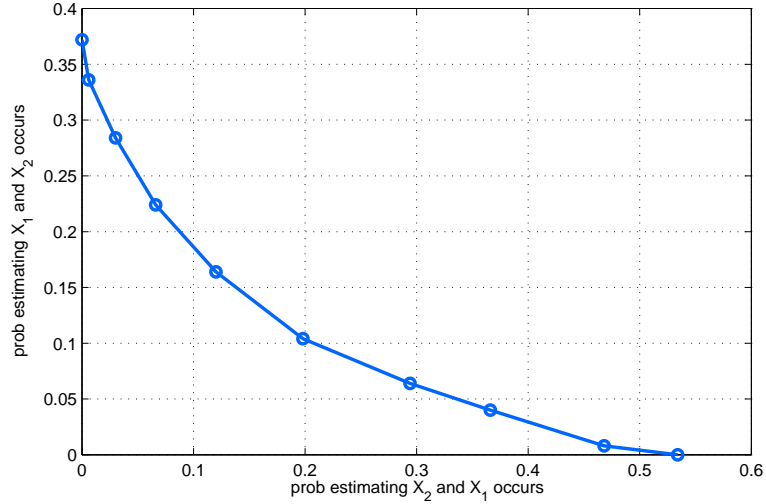
The probability of a false positive is  $E_{21} = 0.366$ . The probability of a false negative is  $E_{12} = 0.040$  Probability of a correct result is  $\text{trace}(E) = 0.5940$ .

The other classifier estimates  $X_1$  if 1, 2, 3, 4, 5 pings are returned, else it estimates  $X_2$ . The error matrix is

$$E = \begin{bmatrix} 0.306 & 0.064 \\ 0.294 & 0.336 \end{bmatrix}$$

The probability of a false positive is  $E_{21} = 0.294$ . The probability of a false negative is  $E_{12} = 0.064$ . Probability of a correct result is  $\text{trace}(E) = 0.642$ .

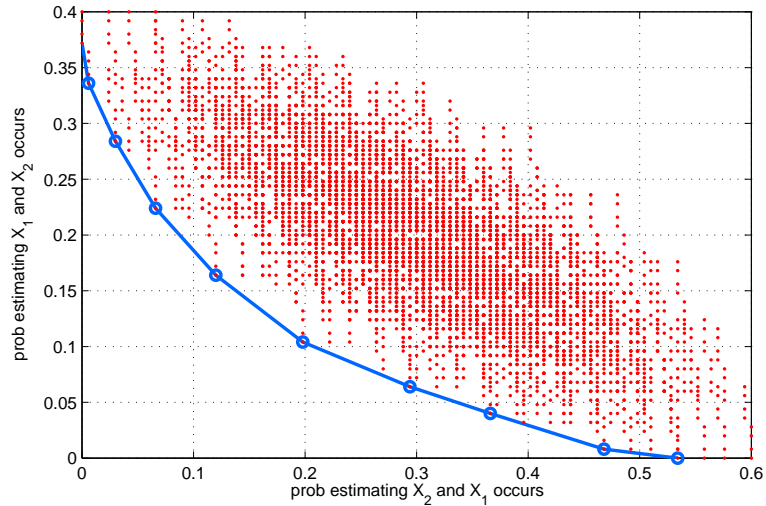
(c) The operating characteristic is below.



(d) The smallest  $\mu$  that achieves  $J_2 < 0.1$  is approximately  $\mu_{\text{opt}} \approx 2.5$ . The classifier is a threshold rule that estimates  $X_1$  for less than 6 pings and  $X_2$  otherwise. The probability of a false positive is  $E_{21} = 0.294$ . The probability of a false negative is  $E_{12} = 0.064$ . Probability of a correct result is  $\text{trace}(E) = 0.642$ .

(e) For each measured number of pings, we can estimate either  $X_1$  or  $X_2$ . Since the number of pings is an element of  $\{1, 2, \dots, 13\}$ , we have  $2^{13}$  different classifiers.

(f) The operating characteristic together with all possible classifiers is below.



(g) For every classifier  $K$ , there exists another classifier  $\hat{K}$  given by

$$\hat{K}_{ij} = 1 - K_{ij}$$

That is, the classifier  $\hat{K}$  decides  $X_2$  whenever  $K$  decides  $X_1$ , and vice versa.

Since  $E = K^T J$ , the unconditional matrix of  $\hat{K}$  is

$$\hat{E} = \begin{bmatrix} \mathbf{Prob}(X_1) - E_{11} & \mathbf{Prob}(X_2) - E_{12} \\ \mathbf{Prob}(X_1) - E_{21} & \mathbf{Prob}(X_2) - E_{22} \end{bmatrix}$$

Hence for each point  $(J_1, J_2)$  on the plot, there is another point at  $(0.6 - J_1, 0.4 - J_2)$ .

- (h) From the previous part, for each estimator with probability of correct result  $p$ , there exists an estimator with probability of error  $p$ . Hence we can look for an estimator which minimizes the probability of error, and then use the opposite classifier to obtain the worst estimator.

We found the best estimator in part(a). From this, the worst estimator is as follows. Estimate  $X_1$  if more than 7 pings are received, else estimate  $X_2$ .

- (i) The optimal decision rule is obtained by computing  $E$  for all the  $2^{13}$  classifiers, and looking for the classifier that minimizes  $J_1$ , subject to  $J_2 < 0.1$ . The optimal classifier is: estimate  $X_1$  if 1, 2, 4, 5, 6 pings are received, else estimate  $X_2$ . For this classifier,  $J_1 = 0.2640$ ,  $J_2 = 0.0960$ , and probability of correct result is 0.64. This classifier was not found in part(d), since the point in the  $(J_1, J_2)$  plane corresponding to this classifier does not lie on the boundary of the convex hull of all possible achievable cost pairs. Hence it cannot be found by minimizing a linear function of  $J_1, J_2$ .
- (j) As the separation in the pmfs increases, the operating characteristic improves, i.e. has a sharper corner, as shown below.

