

# BDA - Assignment 1

## Contents

Exercise 1.. . . . .	1
Exercise 2.. . . . .	1
Exercise 3.. . . . .	3
Exercise 4.. . . . .	3
Exercise 5.. . . . .	4

## Exercise 1.

Probability is the likelihood or chance of an event occurring.

A probability mass function (PMF) gives probabilities for discrete random variables.

So a probability mass is the probability for a discrete random variable.

Probability density is the density of a probability distribution of a continuous variable. Accumulation of probability density gives probability.

Probability density function (PDF) is a statistical expression that defines a probability distribution for a continuous random variable as opposed to a discrete random variable.

A probability distribution is a function that describes the likelihood of obtaining the possible values that a random variable can assume.

A discrete distribution describes the probability of occurrence of each value of a discrete random variable.

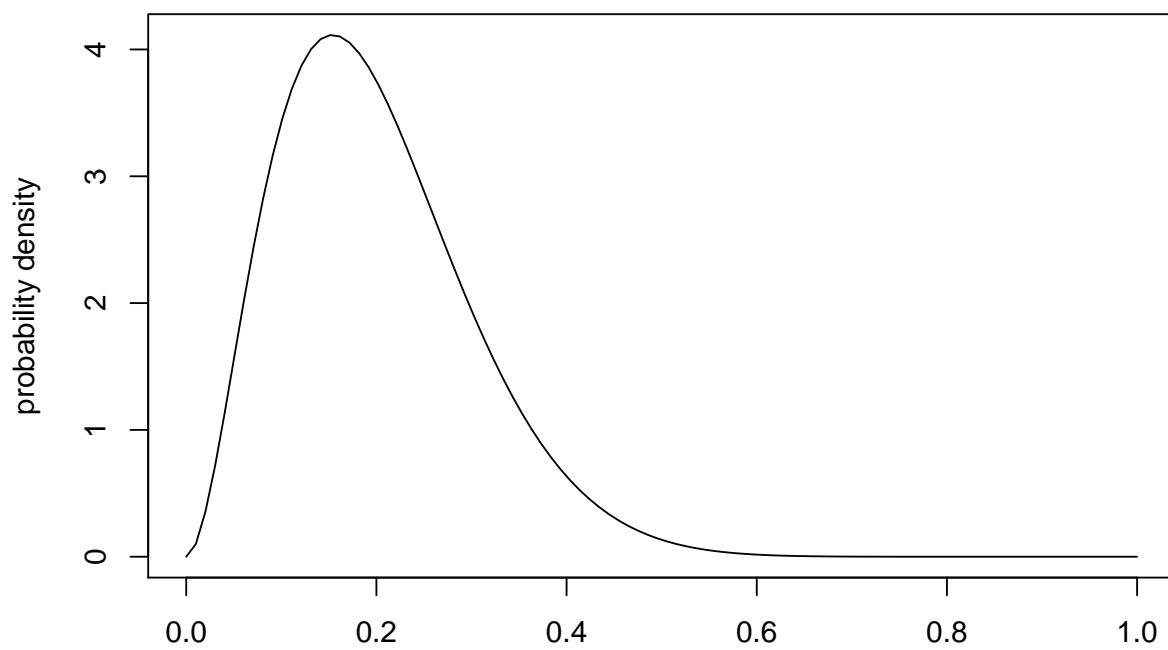
So then Continuous probability distribution is a type of distribution that deals with continuous types of data or random variables.

A likelihood is a function of parameters within the parameter space that how probable a given set of observations is for different values of statistical parameters.

## Exercise 2.

a)

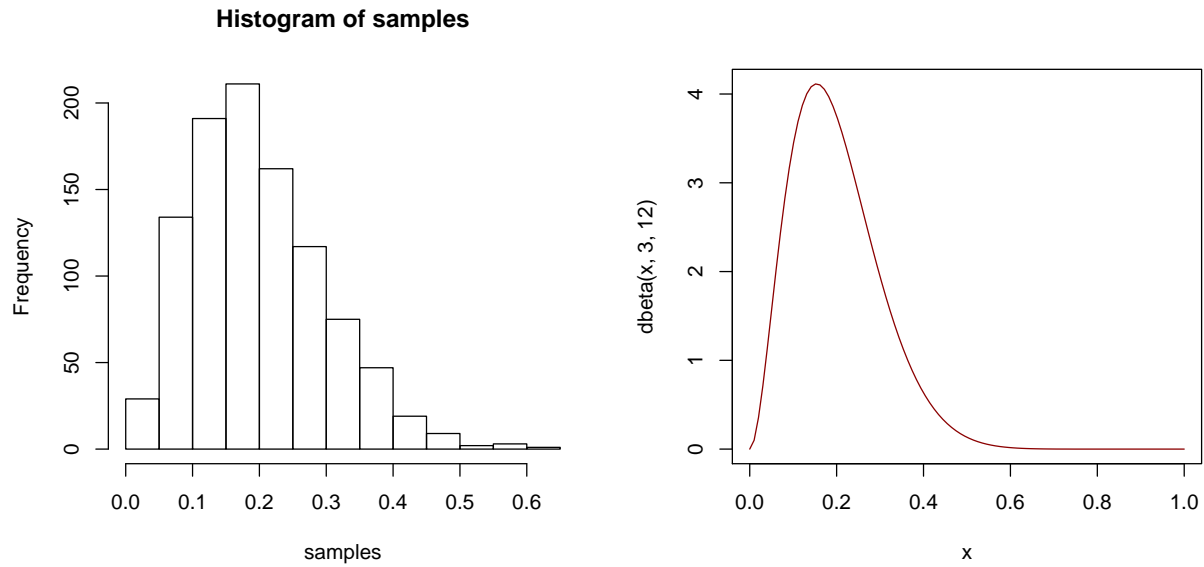
```
x<-seq(0,1,length=100)
plot(x,dbeta(x,3,12),"l",sub="density function of Beta-distribution with
      alpha=3 and beta=12",ylab="probability density")
```



density function of Beta-distribution with  
 $\alpha=3$  and  $\beta=12$

b)

```
par(mfcol=c(1,2))
samples<-rbeta(1000,3,12)
hist(samples)
plot(x,dbeta(x,3,12),"l",col = "dark red")
```



c)

```
mean(samples)
```

```
## [1] 0.1975164
```

```
var(samples)
```

```
## [1] 0.009854737
```

d)

```
quantile(samples, probs=c(0.025, 0.975))
```

```
##      2.5%      97.5%
```

```
## 0.04786022 0.42105764
```

### Exercise 3.

Basically we have several given probabilities:  $P(\text{Positive}|\text{Cancer}) = 0.98$ ,  $P(\text{Negative}|\overline{\text{Cancer}}) = 0.96$  and  $P(\text{Cancer}) = 0.001$ . Taking about success rate, actually if we think over it carefully again, the success rate should be  $P(\text{success}) = P(\text{cancer}|\text{positive}) = \frac{P(\text{positive}|\text{cancer}) * P(\text{cancer})}{P(\text{positive})}$  where  $P(\text{positive}) = P(\text{cancer}) * P(\text{positive}|\text{cancer}) + P(\overline{\text{cancer}}) * (1 - P(\text{negative}|\overline{\text{cancer}})) = 0.04094$ . So then  $P(\text{success}) \approx 0.024 = 2.4\%$  which is much lower than what the researchers expected. Based on the formula, my suggestion is that they should lower down the rate of false positive, it is because that in most cases the examinees don't have cancer so decrease the number of false positive will increase the success rate much more than increasing the ability to distinguish cancer patients from examinees.

### Exercise 4.

Formula used in a):  $P(\text{red}) = P(A) * P(\text{red}|A) + P(B) * P(\text{red}|B) + P(C) * P(\text{red}|C)$

Formula used in b):  $P(A|\text{red}) = \frac{P(\text{red}|A) * P(A)}{P(\text{red})}$   $P(B|\text{red}) = \frac{P(\text{red}|B) * P(B)}{P(\text{red})}$   $P(C|\text{red}) = \frac{P(\text{red}|C) * P(C)}{P(\text{red})}$

```
boxes <- matrix(c(2,4,1,5,1,3), ncol = 2,
dimnames = list(c("A", "B", "C"), c("red", "white")))
```

```
p_red<-function(boxes){
pred1=boxes[1,1]/(boxes[1,1]+boxes[1,2])
pred2=boxes[2,1]/(boxes[2,1]+boxes[2,2])
pred3=boxes[3,1]/(boxes[3,1]+boxes[3,2])
return(pred1*0.4+pred2*0.1+pred3*0.5)
}
```

```
p_box<-function(boxes){
pred1=boxes[1,1]/(boxes[1,1]+boxes[1,2])
pred2=boxes[2,1]/(boxes[2,1]+boxes[2,2])
pred3=boxes[3,1]/(boxes[3,1]+boxes[3,2])
pred=pred1*0.4+pred2*0.1+pred3*0.5
pbox1=(0.4*boxes[1,1]/(boxes[1,1]+boxes[1,2]))/pred
pbox2=(0.1*boxes[2,1]/(boxes[2,1]+boxes[2,2]))/pred
pbox3=(0.5*boxes[3,1]/(boxes[3,1]+boxes[3,2]))/pred
return(c(pbox1,pbox2,pbox3))
}
```

```
p_red(boxes)
```

```
## [1] 0.3192857
```

```
p_box(boxes)
```

```
## [1] 0.3579418 0.2505593 0.3914989
```

## Exercise 5.

Formula used:  $P(\text{identical twins}|\text{twin boys}) = \frac{P(\text{twin boys}|\text{identical twins}) * P(\text{identical twins})}{P(\text{twin boys})}$

$$P(\text{twin boys}) = 0.5 * P(\text{identical twins}) + 0.25 * P(\text{fraternal twins})$$

```
p_identical_twin<-function(fraternal_prob, identical_prob){
return((identical_prob*0.5)/(identical_prob*0.5+fraternal_prob*0.25))
}
p_identical_twin(1/150, 1/400)
```

```
## [1] 0.4285714
```