

# Bayesian Data Analysis - Assignment 1

September 16, 2018

## 1. Basic probability theory notation and terms

a) Explanations:

- *probability* is the estimation of the possibility that an event will occur.
- *probability mass* refers to the probability of samples on an interval; eg. the entire probability sample space is equal to 1.
- *probability density* is the probability of mass divided by unit of the sample space.
- *probability mass function (pmf)* gives the probabilities of the possible values for a discrete random variable.
- *probability density function (pdf)* is a function of a continuous random variable, whose integral across an interval gives the probability that the value of the variable lies within the same interval.
- *probability distribution* is a function of a discrete variable whose integral over any interval is the probability that the variate specified by it will lie within that interval.
- *discrete probability distribution* refers to the probability of occurrence of each value of a discrete random variable.
- *continuous probability distribution* refers to the probabilities of the possible values of a continuous random variable.
- *cumulative distribution function (cdf)* calculates the cumulative probability for a given x-value and it can be used to determine the probability that a random observation that is taken from the sample space will be less than or equal to a certain value.
- *likelihood* is a function of the parameters of a statistical model, given specific observed data.

b) Answers to the questions:

- *What is observation model?*  
Observation model refers to the expression that relates the parameters of the model to the observations.
- *What is statistical model?*  
Statistical modeling is a mathematically simplified way to approximate reality and optionally to make predictions from this approximation.

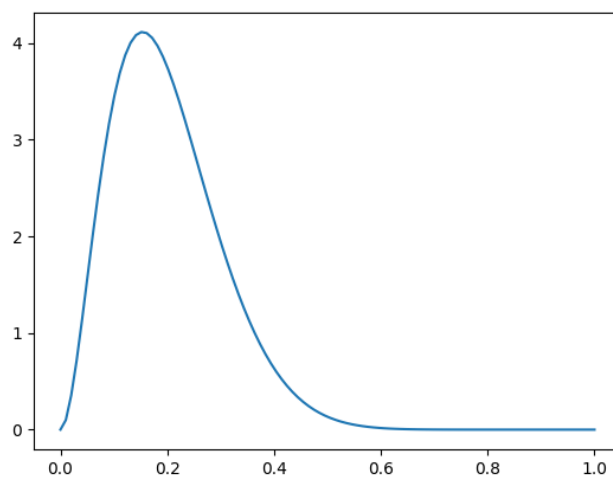
- What is the difference between mass and density?

Mass refers to the entire sample space, whereas density is the fraction from that sample space.

## 2. Basic computer skills

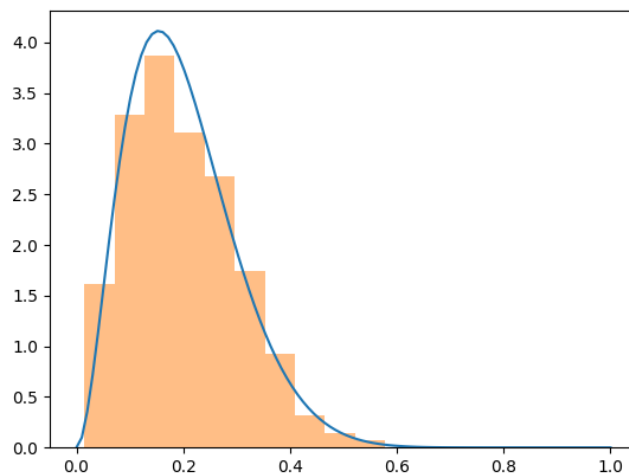
- a) Plot the density function

```
1 from scipy import stats
2 import numpy
3 import matplotlib
4 import matplotlib.pyplot as plt
5
6 MEAN = 0.2
7 VARIANCE = 0.01
8
9 fig, axis = plt.subplots(1, 1)
10
11 alfa = MEAN * ( (MEAN * (1 - MEAN) / VARIANCE) - 1 )
12 beta = alfa * (1 - MEAN) / MEAN
13
14 x_range = numpy.linspace(0, 1, 100)
15 y_range = stats.beta.pdf(x_range, alfa, beta)
16
17 axis.plot(x_range, y_range)
```



- b) Take a sample of 1000 random numbers and plot a histogram of the results

```
1 random_samples = stats.beta.rvs(alfa, beta, size=1000)
2 axis.hist(random_samples, density=True, alpha=0.5)
```



c) Compute the sample mean and variance from the drawn sample

```
1 sample_mean = numpy.mean(random_samples)
2 sample_variance = numpy.var(random_samples)
3 print('sample mean: ', sample_mean)
4 print('sample variance: ', sample_variance)
```

```
$ sample mean:  0.19997418955672838
$ sample variance:  0.01045597802573812
```

d) Estimate the central 95%-interval of the distribution

```
1 sample_percentile = numpy.percentile(random_samples, q=97.5)
2 print('sample central percentile 95%: ', sample_percentile)
```

```
$ sample central percentile 95%:  0.4188436088624379
```

3. **Bayes' theorem:** How would you advice a group of researchers who designed a new test for detecting lung cancer?

$$P(Person_{randomly\ selected\ and\ has\ cancer} | Positive\ result) = \frac{P(P_{has\ cancer}) * P(Cancer)}{P(Positive\ result)} = \frac{0.98 * 0.001}{0.98 * 0.001 + 0.04 * 0.999} = 0.023$$

The result of test is not very satisfying for a randomly selected person who has cancer. Thus, I would tell the researchers to make their test prediction better.

4. **Bayes' theorem:** Find the probability of the selected ball being red and the box it came from.

$A = 2_{red} \ 5_{white}$ ; selected 40% of the time

$B = 4_{red} \ 1_{white}$ ; selected 10% of the time

$C = 1_{red} \ 3_{white}$ ; selected 50% of the time

$$P(\text{red}) = 0.4 * \frac{2}{7} + 0.1 * \frac{4}{5} + 0.5 * \frac{1}{4} = 0.319$$

$$P(A|\text{red}) = \frac{P(\text{red from A}) * P(A)}{P(\text{red})} = \frac{\frac{2}{7} * 0.4}{0.319} = 0.358$$

$$P(B|\text{red}) = \frac{P(\text{red from B}) * P(B)}{P(\text{red})} = \frac{\frac{4}{5} * 0.1}{0.319} = 0.250$$

$$P(C|\text{red}) = \frac{P(\text{red from C}) * P(C)}{P(\text{red})} = \frac{\frac{1}{4} * 0.5}{0.319} = 0.391$$

The probability of a red ball being picked is 31.9% and there is 39.1% chance that it came from Box C.

5. **Bayes' theorem:** What is the probability that Elvis was an identical twin?

Let's assume that:  $S_{gt}$  is a notation of same gender twins,  $I_t$  stands for identical twins and  $F_t$  means fraternal twins. We need to find  $P(I_t|S_{gt})$ .

$$P(S_{gt}) = P(I_t) + P(F_t \text{ for the same gender}) = \frac{1}{300} + 0.5 * \frac{1}{125} = 0.0073$$

$$P(I_t|S_{gt}) = \frac{P(\text{Elvis being twin}) * P(I_t)}{S_{gt}} = \frac{1 * \frac{1}{300}}{0.0073} = 0.45$$

There is 45% chance that Elvis was an identical twin.

## Appendix A Source code

```

1  from scipy import stats
2  import numpy
3  import matplotlib
4  matplotlib.use('TkAgg')
5  import matplotlib.pyplot as plt
6
7  MEAN = 0.2
8  VARIANCE = 0.01
9  fig, axis = plt.subplots(1, 1)
10
11  alfa = MEAN * ( (MEAN * (1 - MEAN) / VARIANCE) - 1 )
12  beta = alfa * (1 - MEAN) / MEAN
13
14  x_range = numpy.linspace(0, 1, 100)
15  y_range = stats.beta.pdf(x_range, alfa, beta)
16
17  # a) Plot the density function of Beta-distribution
18  axis.plot(x_range, y_range)
19  fig.savefig('./ex1/prob_distribution.png')
20
21  # b) Take a sample of 1000 random numbers and plot a histogram
22  random_samples = stats.beta.rvs(alfa, beta, size=1000)
23  axis.hist(random_samples, density=True, alpha=0.5)
24  fig.savefig('./ex1/prob_distribution_hist.png')
25
26  # c) Compute the sample mean and variance from the drawn sample
27  sample_mean = numpy.mean(random_samples)
28  sample_variance = numpy.var(random_samples)
29  print('sample mean: ', sample_mean)
30  print('sample variance: ', sample_variance)
31

```

```
32 # d) Estimate the central 95%-interval from the drawn samples
33 sample_percentile = numpy.percentile(random_samples, q=97.5)
34 print('sample central percentile 95%', sample_percentile)
```