

DATA 5600: Introduction to Regression and Machine Learning for Analytics

Some Brief Notes on Basic Probability Concepts

Author: Tyler J. Brough Updated: October 4, 2021

```
import numpy as np
from scipy import stats
import matplotlib.pyplot as plt

plt.rcParams['figure.figsize'] = [10, 8]
```

Probability Distributions

These notes are based upon readings from the following books:

- *Doing Bayesian Data Analysis* by John Kruschke.
- *Mathematical Statistics with Applications* by Wackerly, Mendenhall, Scheaffer
- *Introduction to Probability and Mathematical Statistics* by Bain & Engelhardt
- *The Cambridge Dictionary of Statistics, 4th Edition*

A **probability distribution** is the list of all possible outcomes and their corresponding probabilities.

Notes: * In class I said that this was a “*mapping from the event space to the probability space*”

- The distribution can be represented by a graph, a table, or a formula.
- Sometimes a distinction is made between the probability *density* and the probability *distribution*, the latter being when the random variable falls at or below some particular value.
- We will use the terms interchangeably and explicitly refer to the latter as the **cumulative distribution function** or simply the **CDF**.

The Discrete Probability Density Function (Probability Mass Function)

If the set of all possible values of a random variable, X , is a countable set, x_1, x_2, \dots, x_n , or x_1, x_2, \dots , then X is called a **discrete random variable**. The function

$$f(x) = P[X = x] \quad x = x_1, x_2, \dots$$

that assigns the probability to each possible value of x will be called the **discrete probability density function** (discrete PDF).

Some Common Examples

Example 1 The values of the discrete pdf of a roll of a fair die can be given by the following table.

x	1	2	3	4	5	6
f(x)	1/6	1/6	1/6	1/6	1/6	1/6

Example 2 When tossing a coin with unknown probability of heads (success). The pdf is given by the ***Bernoulli distribution function***.

$$f(x; \theta) = \theta^x (1 - \theta)^{1-x} \quad x = \{0, 1\}$$

Note:

- $P(X = 0) = 1 - \theta$
- $P(X = 1) = \theta$

A jar contains 30 green jelly beans and 20 purple jelly beans. What is the probability of drawing a single green? A single purple?

$$P(\text{drawing a single green}) = 30/50 = 0.6$$

$$P(\text{drawing a single purple}) = 20/50 = 0.4$$

We can confirm this in **Python** as follows:

```
stats.binom(1, 0.6).pmf(1)
stats.binom(1, 0.6).pmf(0)
```

Example 3 When tossing a coin n times and counting the number of heads the pdf is given by the ***Binomial distribution function***.

$$f(x; n, \theta) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}$$

A jar contains 30 green jelly beans and 20 purple jelly beans. Suppose 10 jelly beans are selected at random from the jar. Find the probability of obtaining exactly five purple jelly beans if they are selected with replacement.

```
stats.binom(10, 0.4).pmf(5)
```

Note: often the definitions of ‘success’ and ‘failure’ are arbitrary. In this case they should be symmetric.

The Cumulative Distribution Function (CDF)

The **cumulative distribution function** (CDF) of a random variable X is defined for any real x by

$$F(x) = P[X \leq x]$$

Example 4 A jar contains 30 green jelly beans and 20 purple jelly beans. Suppose 10 jelly beans are selected at random from the jar. What is the probability of getting 4 or fewer green jelly beans?

```
stats.binom(10, 0.6).cdf(4)
```

Check this against the following and make sure it makes sense to you based on the definitions of a discrete random variable and the CDF given above:

```
p = 0
for i in range(5):
    p += stats.binom(10, 0.6).pmf(i)

print(f"{p : 0.10f}")
```

Q: did you get the same answer? Why or why not?

The Continuous Probability Density Function

A random variable X is called a **continuous random variable** if there is a function $f(x)$ called the **probability density function** (pdf) of X , such that the CDF can be represented as

$$F(X) = \int_{-\infty}^x f(t)dt$$

Example 5 Suppose the research department of a steel manufacturer believes that one of the company's rolling machines is producing sheets of steel of varying thickness. The thickness is represented by a random variable following a uniform random variable with values between 150 and 200 millimeters. Any sheets less than 160 millimeters must be scrapped because they are unacceptable to buyers.

- **a.** Calculate and interpret the mean and standard deviation of x , the thickness of the sheets produced by this machine.
- **b.** Graph the probability distribution of x , and show the mean on the horizontal axis. Also show the 1- and 2-standard deviation intervals around the mean.
- **c.** Calculate the fraction of steel sheets produced by this machine that have to be scrapped.

Solution

a.

To calculate the mean and standard deviation for x , we substitute 150 and 200 millimeters for c and d , respectively, in the formulas for uniform random variables. Thus,

$$\mu = \frac{c + d}{2} = \frac{150 + 200}{2} = 175 \text{ millimeters}$$

and

$$\sigma = \frac{d - c}{\sqrt{12}} = \frac{200 - 150}{\sqrt{12}} = \frac{50}{3.464} = 14.43 \text{ millimeters}$$

b.

```
x = np.linspace(150, 200, 1000)
y = stats.uniform.pdf(x, loc=150, scale=50)
plt.plot(x, y, lw = 2.0, color='darkblue', alpha=0.8)
plt.fill_between(x, y, facecolor='orange', alpha=0.5)
plt.title(f"The Uniform Distribution")
plt.show()
```

c.

To find the fraction of steel sheets produced by the machine that have to be scrapped, we must find the probability that x , the thickness, is less than 160 millimeters. We need to calculate the area under the frequency function $f(x)$ points $x = 150$ and $x = 160$. Therefore, in this case $a = 150$ and $b = 160$.

We have

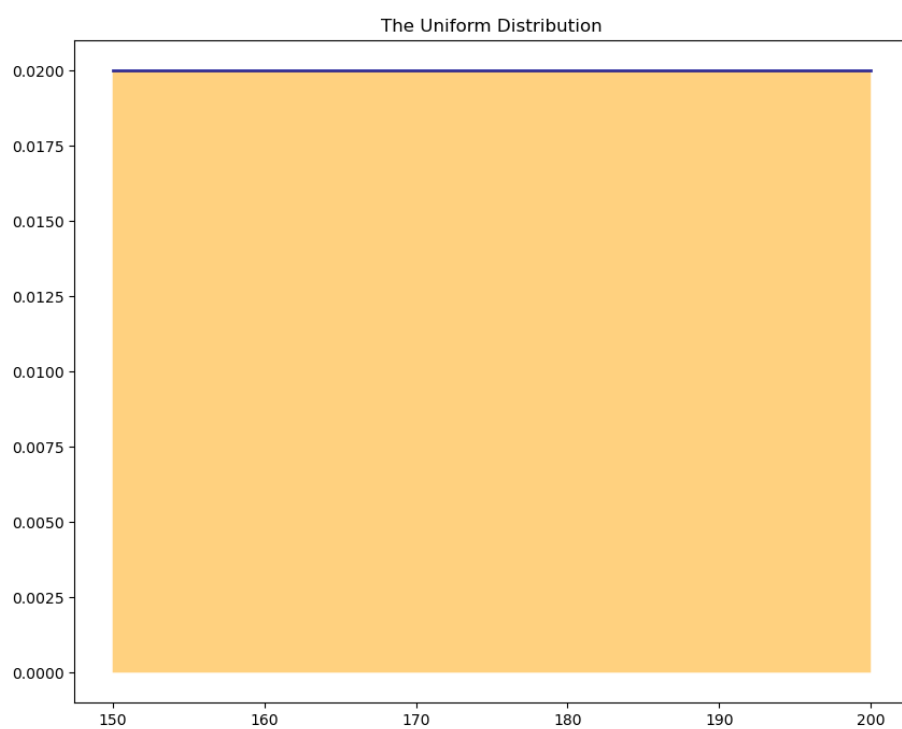


Figure 1: png

$$P(x < 160) = P(150 < x < 160) \\ = \frac{b-a}{d-c} = \frac{160-150}{200-150} = \frac{10}{50} = \frac{1}{5} = 0.20$$

That is, 20% of all the sheets made by this machine must be scrapped.

Of course, we can also simply use the Uniform CDF.

```
start = 150
width = 50
stats.uniform(loc=start, scale=width).cdf(160)
```

Example 6 Suppose the length of time (in hours) between emergency arrivals at a certain hospital is modeled as an *exponential distribution* with $\lambda = 2$. What is the probability that more than 5 hours pass without an emergency arrival?

Solution

The probability we want is the area under the curve to right of 5. To find this probability we use the CDF function and the complement rule.

```
P(X > 5) = 1 - P(X ≤ 5) = 1 - .917915 = .082085

## Set the rate parameter
lam = 2.0

## Use the complement rule
1 - stats.expon(scale=lam).cdf(5)
```

Example 7

The Central Limit Theorem

Let X_1, X_2, \dots, X_n be independent and identically distributed (*iid*) random variables with $E(X_i) = \mu$ and $V(X_i) = \sigma^2 < \infty$. Define

$$Z_n = \sqrt{n} \left(\frac{\bar{X} - \mu}{\sigma} \right) \quad \text{where} \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

Then the distribution function of Z_n converges to a standard normal distribution as $n \rightarrow \infty$.

That is, $Z_n \xrightarrow{d} Z \sim N(0, 1)$ as $n \rightarrow \infty$.

Example n